

Naïve Bayes and Support Vector Machine



Let us understand this concept with the help of an example.

Suppose we have a data set of customers. In this dataset age and income are the attributes and buys_computer as a class with values as Yes and No.

Let us consider a sample X | age : 45 years, income : Rs. 50,000. we have to classify this unseen tuple as buys_computer = yes or buys_computer = No.

$P(H)$ is the probability that the customer X will buy a computer given his age and income, in our case age = 45 and income = Rs. 50,000.

$P(H)$ is the probability any customer will buy a computer regardless of his age and income, it is independent of X . $P(X)$ is the probability from the data set that the customer is 45 years old and earns Rs. 50,000.

$P(X|H)$ is the probability that a customer is 45 years old and he earns Rs. 50,000, given that he will buy a computer.

4.1.2 Maximum A Posteriori (MAP) Hypothesis

Based on Bayes Rule, we can compute the Maximum A Posteriori Hypothesis for the data:

$$h_{\text{map}} = \arg \max_{h \in H} P(h|D)$$

$$= \arg \max_{h \in H} \frac{P(D|h)P(h)}{P(D)} = \arg \max_{h \in H} P(D|h)P(h)$$

H : set of all Hypothesis

4.1.3 Maximum Likelihood

Now assume that all hypotheses are equally probable a priori, i.e., $P(h) = P(h_j)$ for all $h, h_j \in H$.

This is called assuming a uniform prior. It simplifies computing the posterior.

$$h_{\text{ml}} = \arg \max_{h \in H} P(D|h)$$

This hypothesis is called maximum likelihood hypothesis.

Syllabus Topic : Naïve Bayes' Classifiers

4.1 Introduction to Naïve Bayes

- Naïve Bayes is a classification algorithm based on Bayes Theorem and the Maximum A Posterior (MAP) hypothesis.
- It is useful for large data sets.
- It is also called as Idiot Bayes and Simple Bayes.
- Naïve Bayes makes an assumption that the effect of an attribute value on a given class is independent of the values of the other attributes. This assumption is known as class conditional independence.
- The above assumption is made to simplify the computation and in this sense it's called as Naïve (it means the data attributes are independent).

Syllabus Topic : Bayes' Theorem

4.1.1 Bayes Theorem

- Consider $X = [x_1, x_2, x_3, \dots, x_n]$ as a sample, whose components represent values made on a set of n attributes.
- In Bayes Theorem X is considered evidence. Let H be a Hypothesis, such that data X belongs to a class C .
- In classification, the goal is to determine $P(H|X)$, the probability that the hypothesis H holds given the evidence (observed data sample X).
- In other words, the probability that sample X belongs to class C .
- $P(H|X)$ is called as a Posterior probability of H conditioned on X , and is given as :

$$P(H|X) = \frac{P(X|H)P(H)}{P(X)}$$

Where $P(H)$ is the a priori Probability of H .

$P(X|H)$ is the a posteriori probability of X conditioned on H .

- Following is the working of Naïve Bayesian Classifier :
1. Let D be a training set of tuples along with their class labels. Each tuple is represented by an n -dimensional vector, $X = (x_1, x_2, x_3, \dots, x_n)$, consisting of n measurements on the tuple in n attributes A_1, A_2, \dots, A_n .
 2. Consider there are m classes, C_1, C_2, \dots, C_m . Given a tuple X , the classifier will predict that X belongs to a class having highest posterior probability conditioned on X which is given as

$$P(C_i|X) > P(C_j|X) \quad \text{for } 1 \leq j \leq m, j \neq i$$

Hence we maximize $P(C_i|X)$. The class C_i for which $P(C_i|X)$ is maximized is called the maximum posterior hypothesis.

By Bayes theorem we have

$$P(C_i|X) = \frac{P(X|C_i)P(C_i)}{P(X)}$$

3. As $P(X)$ is constant for all classes, only the $P(X|C_i)P(C_i)$ needs to be maximized. In case the prior probabilities of class are not known then it is likely that the probabilities are equal i.e. $P(C_1) = P(C_2) = P(C_3) = \dots = P(C_m)$ and therefore only $P(X|C_i)P(C_i)$ needs to be maximized else $P(X|C_i)P(C_i)$ needs to be maximized. The class prior probabilities can be calculated as

$$P(C_i) = |C_{i,D}| / |D|$$

Where $|C_{i,D}|$ is the number of training tuples of class C_i in D .

4. It is computationally expensive to compute $P(X|C_i)$ given a data set with many attributes. To reduce the computational complexity, naïve assumption of class conditional independence is made (No dependence relationships among attributes), thus

$$P(X|C_i) = \prod_{k=1}^n P(x_k|C_i)$$

$$= P(x_1|C_i) \times P(x_2|C_i) \times \dots \times P(x_n|C_i)$$

These probabilities can be estimated $P(x_1|C_i)$, $P(x_2|C_i)$, $P(x_3|C_i)$, ..., $P(x_n|C_i)$ from the training tuples. x_k is the value of attribute A_k for tuple X .

Computation of $P(x_k|C_i)$ is dependent on whether the attribute is categorical or continuous – valued.

- (a) If A_k is Categorical : $P(x_k|C_i)$ is the number of tuples of class C_i in D having the value x_k for A_k , divided by $|C_{i,D}|$, the number of tuples of class C_i in D .

- (b) If A_k is Continuous : A continuous valued attribute has a Gaussian distribution with a mean μ and standard deviation σ , defined by

$$g(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$P(x_k|C_i) = g(x_k, \mu_{C_i}, \sigma_{C_i})$$

Compute μ_{C_i} and σ_{C_i} , which are the mean (i.e., average) and standard deviation, respectively, of the values of attribute A_k for training tuples of class C_i .

- For Prediction of class label of X , $P(X|C_j)P(C_j)$ is evaluated for each class C_j . The classifier predicts that the class label of tuple X is the class C_j if and only if

$$P(X|C_j)P(C_j) > P(X|C_i)P(C_i) \text{ for } 1 \leq j \leq m, j \neq i.$$

Which means the predicted class label is the class C_j for which $P(X|C_j)P(C_j)$ is the maximum.

4.2.1 E.g. Naïve Bayesian Classifier

Predict a class label of an unknown sample using Naïve Bayesian classification on the following training dataset from all electronics customer database.

Given a training set, we can compute the probabilities as follows :

The classification problem may be formalized using a posterior probabilities :

$P(C|Y)$ is the probability that the sample tuple $Y = \langle y_1, \dots, y_n \rangle$ is of class C .

- Assign to sample Y the class label C such that $P(C|Y)$ is maximal.

Age	Income	Student	Credit rating	Class buys computer
≤ 30	High	No	Fair	No
≤ 30	High	No	Excellent	No
$31 \dots 40$	High	No	Fair	Yes
> 40	Medium	No	Fair	Yes
≥ 40	Low	Yes	Fair	Yes
≥ 40	Low	Yes	Excellent	No
$31 \dots 40$	Low	Yes	Excellent	Yes
≤ 30	Medium	No	Fair	No
≤ 30	Low	Yes	Fair	Yes
> 40	Medium	Yes	Fair	Yes
≤ 30	Medium	Yes	Excellent	Yes
$31 \dots 40$	Medium	No	Excellent	Yes
$31 \dots 40$	High	Yes	Fair	Yes
> 40	Medium	No	Excellent	No

The unknown sample is $x = \text{age} = \text{"}\leq 30\text{"}$, income = "Medium", student = "Yes", credit rating = "Fair"

Age
≤ 30 Yes
≤ 30 No
$31 \dots 40$ No
$31 \dots 40$ Yes
> 40 No
> 40 Yes
≤ 30 No
≤ 30 Yes
$31 \dots 40$ No
$31 \dots 40$ Yes
> 40 No
> 40 Yes

$P(\leq 30 \text{Yes}) = 29$	$P(\leq 30 \text{No}) = 35$
$P(31 \dots 40 \text{Yes}) = 49$	$P(31 \dots 40 \text{No}) = 0$
$P(>40 \text{Yes}) = 39$	$P(>40 \text{No}) = 25$
Income	
High	$P(\text{High} \text{Yes}) = 25$
Medium	$P(\text{Medium} \text{Yes}) = 49$
Low	$P(\text{Low} \text{Yes}) = 39$
	$P(\text{Low} \text{No}) = 15$

Machine Learning (SPPU - Sem 8 - Comp.)		4.5	Naïve Bayes and Support Vector Machine
Student			

P(Yes yes) = 3/9	P(No no) = 4/5
P(Yes no) = 6/9	P(Yes no) = 1/5
Credit Rating	
P(fair yes) = 6/9	P(fair no) = 2/5
P(excellent yes) = 3/9	P(excellent no) = 3/5
P(fair no) = 9/14	
P(excellent no) = 5/14	

An unseen sample $X = \langle \text{age} = <30, \text{income} = \text{Medium}, \text{Student} = \text{Yes}, \text{Credit rating} = \text{Fair} \rangle$

$$P(X|Y_{\text{Yes}}) \cdot P(Y_{\text{Yes}}) = P(\text{Age} = <30 | \text{Yes}) \cdot P(\text{income} = \text{Medium} | \text{Yes})$$

$$\cdot P(\text{Student} = \text{Yes} | \text{Yes}) \cdot P(\text{Credit Rating} = \text{Fair} | \text{Yes})$$

$$P(\text{Yes}) = 29/49 \cdot 6/9 \cdot 6/9 \cdot 9/14 = 0.028$$

$$P(X|Y_{\text{No}}) \cdot P(Y_{\text{No}}) = P(\text{Age} = <30 | \text{No}) \cdot P(\text{income} = \text{Medium} | \text{No})$$

$$P(\text{No}) = 0.007$$

Since $0.028 > 0.007$, Therefore the naive Bayesian classifier predicts Buys computer = "Yes" for sample X

4.2.2 Properties of Bayes Classifiers

1. Instrumentality

- The prior and likelihood can be updated dynamically with each training example.

2. Combines prior knowledge and observed data

- Given the training data, the prior probability of a hypothesis is multiplied with probability of the hypothesis.

3. Probabilistic hypotheses

- The output of Bayes classifier includes classification as well as probability distribution over all classes.

4. Meta-classification

- The output of several classifiers can be combined together.
- For e.g. the probabilities of all classifiers predicted for a given class can be multiplied together.

Machine Learning (SPPU - Sem 8 - Comp.)		4.6	Naïve Bayes and Support Vector Machine
Syllabus Topic : Naïve Bayes in Scikit-Learn			

Machine Learning (SPPU - Sem 8 - Comp.)		4.6	Naïve Bayes and Support Vector Machine
Syllabus Topic : Bernoulli Naïve Bayes			

4.3 Bayes in Scikit-Learn

- Based on the same number of different probabilistic distributions, scikit - learn provides implementation of three variants of Naïve bayes : Bernoulli Naïve Bayes, Multinomial Naïve Bayes and Gaussian Naïve Bayes.
- Bernoulli naïve bayes is a binary distribution and this type is useful when a feature is present or absent.
- Multinomial Naïve Bayes is a discrete distribution and is useful when a feature needs to represented using a whole number.
- The Gaussian distribution is a continuous distribution and is characterised by mean and variance.

Syllabus Topic : Bernoulli Naïve Bayes

4.4 Bernoulli Naïve Bayes

- Consider X is a random variable having Bernoulli distribution.

It can assume only two values (say for e.g. 0 and 1) and their probability is given as follows :

$$P(X) = \begin{cases} p & \text{if } X = 1 \\ q & \text{if } X = 0 \end{cases}$$

Where $q = 1 - p$ and $0 < p < 1$

E.g. : Implementation of Bernoulli Naïve Bayes in Scikit-Learn

- Generating a dummy data set :

```
from sklearn.datasets import make_classification
n_samples = 300
X, Y = make_classification(n_samples=300, n_features=2, n_informative=2, n_redundant=0)
```

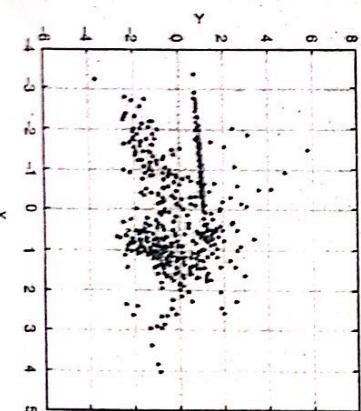


Fig. 4.4.1 : Dummy data generated

Machine Learning (SPPU - Sem 8 - Comp.)

Naïve Bayes and Support Vector Machines

- The `binarize` parameter in `BernoulliNB` class allows using a threshold that can be used internally to transform the features to binary
- Let us use 0.0 as the binary threshold.

```
from sklearn.naive_bayes import BernoulliNB
```

卷之三

#OUTPUT

	0	(2,1)	(1,1)
(0,0)			
(1,0)			

Class 1

indina

Syllabus Topic : Multinomial Naïve Bayes

三

- Gaussian Naive Bayes is used when the values are continuous whose probability distribution

4.5 Multinomial Naïve Bayes

1

Conditional Probabilities are also Gaussian Distributed, mean and variance needed

If a feature vector has n elements and each of them have k different values with probability p_i , then

$$P(X_1 = x_1 \cap X_2 = x_2 \cap \dots \cap X_t = x_t) = \frac{n!}{\prod_i x_i!} \prod_i p_i^{x_i}$$

To prevent the model from setting null probabilities when the frequency is zero, an alpha parameter (called Laplace smoothing factor) is used whose default value is 1.0

The output class for city is 1 and countryside is 0

```

from sklearn import datasets
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn import metrics

# Load the dataset
iris = datasets.load_iris()
X = iris.data[:, [0, 2]] # Features
y = iris.target # Target variable
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.5, random_state=42)
knn = KNeighborsClassifier(n_neighbors=3)
knn.fit(X_train, y_train)
y_pred = knn.predict(X_test)

# Evaluate the model
print("Accuracy: ", metrics.accuracy_score(y_test, y_pred))

```

Machine Learning (SPPU : Sem B : Comp)

AMA Bayas and Support Vector Machines

- Bernoulli naive Bayes needs binary feature vectors.
- The `binarize` parameter in `BernoulliNB` class allows using a threshold that can be used internally to transform the features to binary
- Let us use 0.0 as the binary threshold.

```
from sklearn.naive_bayes import BernoulliNB
```

卷之三

#OUTPUT

卷之三

Scanned with CamScanner

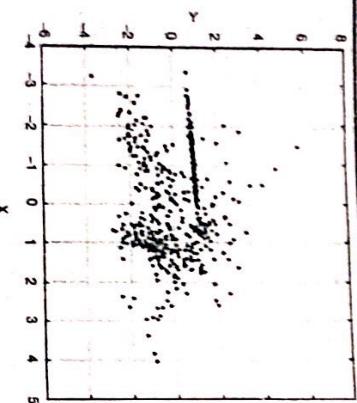


Fig. 4.6.1 : A plot of dataset

Training both the models and generating ROC curves

```
from sklearn.naive_bayes import GaussianNB
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn import metrics
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.25) >>> gnb = GaussianNB()
gnb.fit(X_train, Y_train)
Y_pred_gnb = gnb.predict_proba(X_test)
lrc = LogisticRegression()
lrc.fit(X_train, Y_train)
Y_pred_lrc = lrc.predict_proba(Y_test)
Auc_gnb = metrics.roc_auc_score(Y_test, Y_pred_gnb[:, 1])
Auc_lrc = metrics.roc_auc_score(Y_test, Y_pred_lrc[:, 1])
```

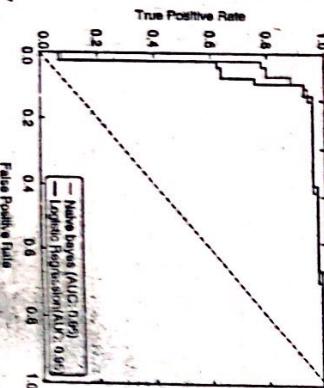


Fig. 4.6.2

As can be seen from the Fig. 4.6.2, Naive bayes performs slightly better than Logistic regression, both the classifiers have similar accuracy and AUC (Area Under Curve).

4.7 Application of Naive Bayes

1. Real-time Prediction

Due to high speed of Naive bayes algorithm, it can be used for making predictions in real time.

2. Multi-class Prediction

The posterior probability of multiple classes of the target variable can be predicted by the algorithm.

3. Text classification/ Spam Filtering/ Sentiment Analysis

- Due to their better results in multiclass problems and independence rule, Naive Bayes classifier is suitable for text classification.

- It has a high success rate compared to other algorithms.

- Due to its success rates they are widely used in spam filtering and sentiment analysis applications.

4. Recommendation System

- Native Bayes can be combined together with collaborative filtering to make a recommendation system in machine learning and data mining to filter the unseen information and predict the users likes and dislikes.

Algorithm Advantages

- The algorithm performs well for multiclass prediction.

- The prediction of class on test data is faster and it is also easy to apply.

- A Naive Bayes classifier performs better when compared to other classifiers like e.g. Logistic regression when the assumption of independence holds, it also requires less amount of training data.

- The algorithm performs better for categorical input variable(s) when compared to numerical variable(s). The numerical variable assumes a normal distribution.

Algorithm Disadvantages

- The model will be unable to make a prediction if the categorical variable has a category in test data, which was not available in training data.

- In such a case the model assigns a zero (0) probability. This is called as "Zero Frequency". This problem can be overcome with a smoothing technique. One of the simplest method used for smoothing is called Laplace estimation.

- The assumption of independent predictors in Naive Bayes algorithm is a limitation as in real world it is almost impossible to get a set of predictors which are independent.

Syllabus Topic : Support Vector Machines (SVM)

4.8 Introduction to Support Vector Machine (SVM)

- Support vector Machine is another method used for classification.

- It can classify both Linear as well as Non Linear Data.

- The objective of SVM is to find a hyperplane (a decision boundary separating the tuples of one class from another) in an N -dimensional space (where N represents the number of features) that distinctly classifies the data points.

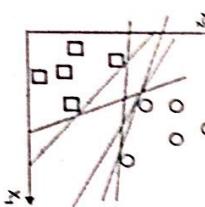


Fig. 4.8.1 : All possible Hyperplanes to classify the data points

- As we can see in the above diagram, there are many possible hyperplanes that could be chosen, but the main objective is to find a plane that has maximum margin i.e. maximum distance between data points of both classes as shown in the Fig. 4.8.2.

- Maximizing Γ provides reinforcement so that future data points can be classified with more confidence.

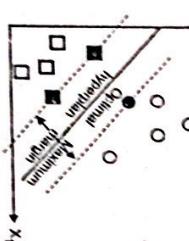


Fig. 4.8.2 : Optimal hyperplane having maximum margin

or Working of Support Vector Machine

- SVM uses a Non Linear mapping to transform the original training data into a higher dimension.

- In this dimension it searches for a linear optimal separating hyperplane.
- It finds this hyperplane using support vectors (essentially training tuples) and margins (defined by support vectors).

- In SVM, if the output of the linear function is greater than 1, it is identified as one class and if the output is -1 it is identified with another class.

or Hyperplanes

- Hyperplanes are decision boundaries that are used to classify the data points.

- The dimension of the hyperplane depends upon the number of features.

- Number of input features is 2 ; then hyperplane is a Line
- Number of input features is 3 ; the hyperplane is a two dimensional plane

- Difficult to imagine a hyperplane when input features are more than 3.

- A hyperplane in R^2 is a line
- A hyperplane in R^3 is a plane

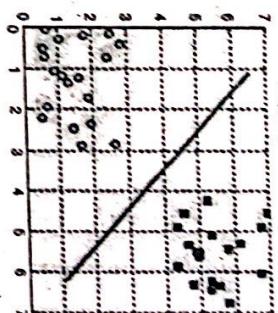


Fig. 4.8.3 : Hyperplane as a Line and a two dimensional plane

or Support Vector

- Support vectors are data points that are closer to hyperplane.
- These support vectors are used to maximize the margin of the classifier.

Deletion of these support vectors will change the position of the Hyperplane.



Fig. 4.8.4 : Support vectors and Margins separating the data points

4.8.1 Tuning Hyperparameters

or Kernel

- Kernel is used to transform the given data into the required form.

- Different types of kernel functions may be used in Support vector machine, such as Linear, Polynomial, Radial Basis function and sigmoid kernel.

- Polynomial and RBF kernels compute the separation line in higher dimension, they are also useful for non linear hyperplane.

- For separating the classes that are either curved or non linear, more complex kernels can be used.

or Regularization

- A penalty parameter called as C is used in Scikit learn library as a regularization parameter.

- The value of C represents misclassification or error term, which tells the SVM optimization by how much the error is bearable. This is how the trade-off between the decision boundary and misclassification term is controlled.
- A small-margin hyperplane is created by a small value of C and a larger value of C creates a larger-margin hyperplane.

Gamma

- The value of gamma decides the number of data points in the calculation of separating line.
- A low value of gamma will consider only nearby points whereas a high value of gamma will consider all the data points in the calculation.
- A low value of gamma will loosely fit the training data, whereas a high value will overfit the training data.

Advantages of SVM

- SVM is a highly accurate method due to the ability of modeling complex non-linear decision boundaries.
- When compared to other methods, SVM is less prone to overfitting.
- The support vectors of the model provide a description of the learned model.
- SVM can be used for prediction (numeric) and also as a classifier.

Disadvantages

- The training time of even the fastest SVMs can be extremely slow on large data sets.
- Less effective on noisier datasets with overlapping classes.

Applications

- Handwritten digit recognition.
- Object recognition.
- Speaker identification.
- Benchmark time-series prediction tests.

Syllabus Topic : Linear Support Vector Machines

4.9 Linear Support Vector Machines

- Let us consider a two class problem where the classes are linearly separable.
- Consider D as the data set which can be given by

$$(X_1, y_1), (X_2, y_2), \dots, (X_m, y_m)$$

Where X_i is the set of training tuples and y_i the associated class labels which can take +1 or -1 as its values.

The Linearly separable classes as shown in the Fig. 4.9.1.

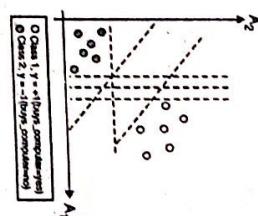


Fig. 4.9.1 : Linearly Separable Data

- Goal here is to find the best possible separating line (considering data as 2D) having minimum classification error on previously unseen tuple.

Incase if the data is 3D, the goal would be to find the best separating hyperplane with minimum classification error.

So to Generalize the concept for N dimensions, the Goal is to find the best Hyperplane.

Support Vector Machine works by finding the maximum marginal hyperplane.

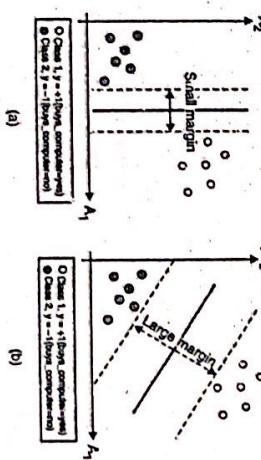


Fig. 4.9.2 : Two possible hyperplanes separating the data

- From Fig. 4.9.2 we can see that there are two possible hyperplanes separating the data, one with smaller margin and the other with larger margin.

The hyperplane with larger margin can classify the data more better in the future as compared to the smaller margin hyperplane.

- This is the reason SVM searches the hyperplane with largest margin, that is Maximal Margin Hyperplane (MMH).
- Definition of Margin : It is the shortest distance from a hyperplane to one side of its margin is equal to the shortest distance from the hyperplane to the other side of its margin, where the "sides" of the margin are parallel to the hyperplane.

- A separating hyperplane can be given as

$$W \cdot X + b = 0$$
- W : is a weight vector, $W = [w_1, w_2, w_3, \dots, w_n]$
- n : is the number of attributes
- b : is a scalar (called as bias)

- To understand the classification let us consider two input attributes A_1 and A_2
- Training tuples are 2D, $(X = (x_1, x_2))$ where x_1 and x_2 are attributes of A_1 and A_2
- Considering w_0 as the additional weight, and rewriting Equation (4.9.1) as

$$w_0 + w_1 x_1 + w_2 x_2 = 0$$

- Any point that lies above the hyperplane satisfies the following equation

$$w_0 + w_1 x_1 + w_2 x_2 > 0$$

- And any point that lies below the hyperplane satisfies the following equation

$$w_0 + w_1 x_1 + w_2 x_2 < 0$$

- Adjusting the weight, and hyperplanes defining the sides of the margin can be written as

$$H_1: w_0 + w_1 x_1 + w_2 x_2 \geq 1 \quad \text{for } Y_1 = +1,$$

$$H_2: w_0 + w_1 x_1 + w_2 x_2 \leq -1 \quad \text{for } Y_1 = -1.$$

From the above two equations (4.9.5) and (4.9.6) we can say that , any tuple that falls on or above H_1 belongs to class +1 and any tuple that falls on or below H_2 belongs to class -1.

Combining the two equations (4.9.5) and (4.9.6) we get

$$Y_1(w_0 + w_1 x_1 + w_2 x_2) \geq 1, \quad \forall i \quad \dots (4.9.7)$$

Any training tuples that fall on hyperplanes H_1 or H_2 (i.e., the "sides" defining the margin) satisfy Eq. (4.9.7) and are called support vectors. That is, they are equally close to the (separating) MMH.

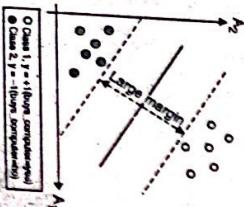


Fig. 4.9.3 Support vector with thick boundary

- The distance from the separating hyperplane to any point on H_1 is $|w|$, where $|w|$ is the Euclidean norm of $\sqrt{w^2}$. By definition, this is equal to the distance from any point on H_2 to the separating hyperplane. Therefore, the maximal margin is $2|w|$.

4.10 Scikit-Learn Implementation

- A support vector machine(SVM) is a type of supervised machine learning classification algorithm used for classification, regression and outlier detection.
- The support vector machines in scikit-learn support both dense and sparse sample vectors as input.
- SVC, NuSVC and LinearSVC are classes capable of performing multi-class classification on a dataset.
- SVC and NuSVC implements "one-against-one" approach for classification , on the other hand LinearSVC implements "one-vs-the-rest" multiclass strategy, thus training n class models. In case of only two classes only one model is trained.

Syllabus Topic : Linear Classification

4.10.1 Linear Based Classification

- E.g. In Scikit-learn, a simple classification task in which two classes of points are well separated.

Standard Imports

```
# Standard Python imports
import numpy as np
import matplotlib.pyplot as plt
```

```
# Import scikit-learn tools
```

```
# Use standard plotting defaults
```

```
# Import scikit-learn as sns
sns.set()
```

```
from sklearn.datasets.samples_generator import make_blobs
```

```
N, d = make_blobs(n_samples=50, centers=2, random_state=0, cluster_std=0.60)
plt.scatter(N[:, 0], N[:, 1], c=s, s=50,
```

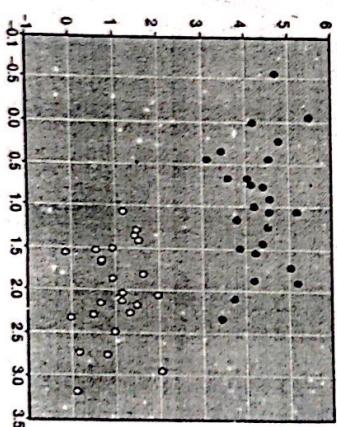


Fig. 4.10.1 : Plot showing data which is linearly separable

Machine Learning (SPPU - Sem 8 - Comp) 4-17

Naïve Bayes and Support Vector Machine

☛ Draw Lines that can divide the data into two separate classes

```
#f1 = np.linspace(-1, 3.5)
P = np.zeros((X.shape[0], Y.shape[0]), dtype='float')
plt.scatter(X[:, 0], X[:, 1], c=Y, s=50, cmap='autumn')
for m, b in [(1, 0.65), (0.5, 1.6), (-0.2, 2.9)]:
    plt.plot(m, b, 'k-')
plt.plot(SG, m * SG + b, 'k-')
```

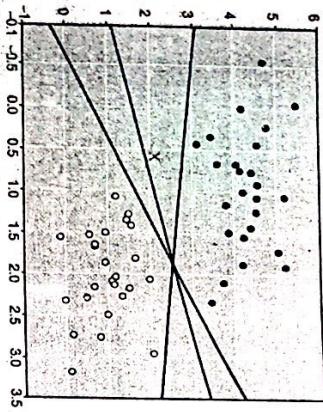


Fig. 4.10.2: Plot showing Lines separating the data

Rather than simply drawing a zero-width line between the classes, we can draw around each line a *margin* of some width, up to the nearest point.

☛ Fitting a support vector machine

Use Scikit-Learn's support vector classifier to train an SVM model on this data.

```
from sklearn import svm
from sklearn import datasets
#Support vector classifier
model = SVC(kernel='linear', C=1E0)
model.fit(X, y)
SVC(C=1000000000.0, dual=False, class_weight=None, degree=3,
gamma='auto', kernel='linear', max_iter=-1, probability=False, random_state=None, shrinking=True,
tol=0.001, verbose=False)
```

(Plot the decision function of model 1, as = None, plot_support = True):

```
#Plot the decision function of model 1, as = None, plot_support = True:
from sklearn import svm
import numpy as np
import matplotlib.pyplot as plt
from sklearn import datasets
iris = datasets.load_iris()
X = iris.data[:, [0, 1, 2, 3]]
y = iris.target
# take the first two features
# we could also do a logistic regression with
# two classes
#np.random.seed(42)
#X = np.random.rand(100, 2)
#y = np.where(np.sum(X, axis=1) < 8, 0, 1)
#np.random.seed(42)

# fit an SVM model not too
# marginally separated
#C = 1.0, epsilon = 0.1
#model = svm.SVC(kernel='linear', C=C, epsilon=epsilon, tol=1e-3)
#model.fit(X, y)

# Plot the decision function of this model
#ax = plt.gca()
#ax.set_title("2-class classification using SVM")
#ax.set_xlim(X[:, 0].min() - 1, X[:, 0].max() + 1)
#ax.set_ylim(X[:, 1].min() - 1, X[:, 1].max() + 1)
#ax.set_xlabel("Sepal length [cm]")
#ax.set_ylabel("Sepal width [cm]")
#ax.set_xticks([0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0])
#ax.set_yticks([0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0])
#ax.grid(True)
#ax.scatter(X[y == 0, 0], X[y == 0, 1], color="red", marker="o", s=50)
#ax.scatter(X[y == 1, 0], X[y == 1, 1], color="blue", marker="x", s=50)
#ax.plot(model.decision_function(X), color="black", linewidth=2)
```

```
#f1 = np.linspace(-1, 3.5)
P = np.zeros((X.shape[0], Y.shape[0]), dtype='float')
plt.scatter(X[:, 0], X[:, 1], c=Y, s=50, cmap='autumn')
for m, b in [(1, 0.65), (0.5, 1.6), (-0.2, 2.9)]:
    plt.plot(m, b, 'k-')
plt.plot(SG, m * SG + b, 'k-')
```

Plot the decision function for 2D SVM:

```
#f1 = np.linspace(-1, 3.5)
P = np.zeros((X.shape[0], Y.shape[0]), dtype='float')
plt.scatter(X[:, 0], X[:, 1], c=Y, s=50, cmap='autumn')
for m, b in [(1, 0.65), (0.5, 1.6), (-0.2, 2.9)]:
    plt.plot(m, b, 'k-')
plt.plot(SG, m * SG + b, 'k-')
```

Machine Learning (SPPU - Sem 8 - Comp) 4-18

Naïve Bayes and Support Vector Machine

☛ Draw Lines that can divide the data into two separate classes

```
#f1 = np.linspace(-1, 3.5)
P = np.zeros((X.shape[0], Y.shape[0]), dtype='float')
plt.scatter(X[:, 0], X[:, 1], c=Y, s=50, cmap='autumn')
for m, b in [(1, 0.65), (0.5, 1.6), (-0.2, 2.9)]:
    plt.plot(m, b, 'k-')
plt.plot(SG, m * SG + b, 'k-')
```

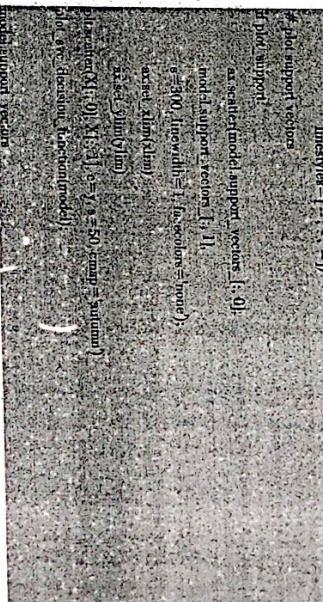


Fig. 4.10.3 : Support vectors

4.10.2 Kernel Based Classification

Syllabus Topic : Kernel Based Classification

- Some problems cannot be solved using Linear hyper-plane for e.g. as shown below.
- In this type of classes, SVM uses a kernel trick to transform the input space to a higher dimensional as shown in the Fig. 4.10.40).

- To find a tradeoff between precision and the number of support vectors, several ways are possible. One way is to implement a parameter α bounded between 0 – not included – and 1) can be used to control the number of support vectors.

卷之三

Fig. 4.10.6 : Linearly separable data using radial kernel

The kernel model hyperparameter in Scikit-learn can be used for applying kernelized SVM, i.e. for changing the Linear kernel to RBF or some other kernel.

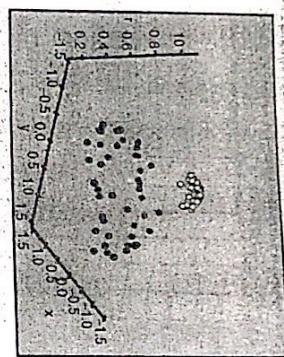


Fig. 4.10.6 : Linearly separable data using Radial Kernel

```

SVC(C=10, kernel='rbf', gamma=0.01, class_weight='auto', degree=3,
      max_iter=-1, probability=False, random_state=None, shrinking=True,
      tol=0.001, verbose=False)

# SVC(X, y, kernel='rbf', degree=3, gamma=0.01, C=10, shrinking=True,
#      tol=0.001, cache_size=100, class_weight=None, probability=False,
#      max_iter=-1, decision_function_shape='ovo', decision_function_shape='ovr',
#      random_state=None, shrinking=True, tol=0.001, verbose=False)

# SVC(C=10, kernel='rbf', gamma=0.01, degree=3, max_iter=-1, probability=False,
#      random_state=None, shrinking=True, tol=0.001, verbose=False)

```

Fig. 4.10.7

Syllabus Topic : Controlled Support Vector Machines

4.11 Controlled Support Vector Machines

While working with real data sets, SVM classifier can result in many support vectors, which can slow down the whole process.

CHAPTER

5

Decision Trees and Ensemble Learning

UNIT V

Syllabus Topics

Decision Trees - Impurity measures, Feature Importances, Decision Tree Classification with Scikit-learn, Ensemble Learning, Random Forest, AdaBoost, Gradient Tree Boosting, Voting Classifier.

Clustering Fundamentals - Basics, K-means, Finding optimal number of clusters, DBSCAN, Spectral Clustering.

Evaluation methods based on Ground Truth- Homogeneity, Completeness, Adjusted Rand Index.

Introduction to Meta Classifier- Concepts of Weak and eager learner, Ensemble methods, Bagging, Boosting, Random Forests.

Syllabus Topic : Decision Trees – Impurity Measures

5.1 Decision Trees

- Suppose a database D is given as $D = \{t_1, t_2, \dots, t_n\}$ and a set of desired classes are $C = \{C_1, \dots, C_m\}$, the classification problem is to define the mapping m in such a way that which tuple of database D belongs to which class of C. Actually divides D into equivalence classes.

Classification Example

How teacher gives grades to students based on their marks obtained:

- If $x > 90$ then grade = A.
- If $80 < x < 90$ then grade = B.
- If $70 < x < 80$ then grade = C.
- If $60 < x < 70$ then grade = D.
- If $x < 60$ then grade = E.

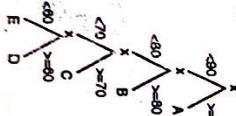


FIG. 5.1.1: Classification of grading

Decision tree learning is appropriate for problems having characteristics given below:

- Instances are represented by a fixed set of attributes (e.g. Gender) and their values (e.g. male, female) described as attribute-value pairs.
- If the attribute has small number of disjoint possible values (e.g. high, medium, low) or there are only two possible classes (e.g. true, false) then decision tree learning is easy.

Machine Learning (SPPU - Sem B - Comb)

5.2

Decision Trees and Ensemble Learning

Extension to decision tree algorithm also handles real valued attributes (e.g. salary).

Decision tree gives a class label to each instance of the dataset.

- If some of the examples have unknown values, still decision tree method can be used (Example: Wind value is known for few tuples).

Decision Tree Representation

Decision tree has leaf nodes and decision nodes.

Every branch's last node is a leaf node which represents the class label.

Decision node is used to take decision for the class label after some trials. It has a leaf node or sub tree.

Play Tennis example for representing decision Trees.

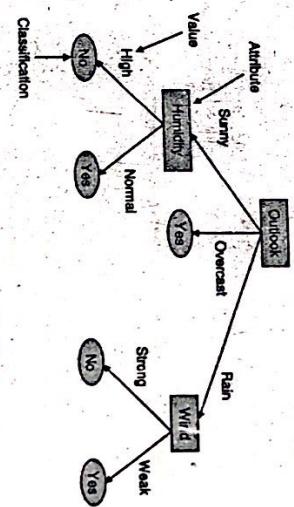


FIG. 5.1.2 : Representation of decision tree

Binary Decision tree

It is a structure based on a sequential decision process.

It goes from the root and based on evaluating features one of the two branches is chosen.

Repeat the same procedure until a final leaf node or classification target is reached.

5.1.1 Impurity measures

1. Gini Impurity Index
 - Suppose all attributes are continuous-valued.
 - Assume that each value of an attribute has many possible splits.
 - It can be adapted for categorical attributes.

- Gini is used in CART (Classification and Regression Trees), IBM's IntelligentMiner system, SPRINT (Scalable Parallelizable Induction of decision Trees).

- If a data set T contains examples from n classes, Gini index, $Gini(T)$ is defined as

$$Gini(T) = 1 - \sum_{j=1}^n p_j^2$$

- Where, b_j is the relative frequency of class j in T .
- $gini(T)$ is minimized if the classes in T are skewed.
- Gini index of split is defined below after split T1 and T2 of sizes N1 and N2 respectively.

$$gini_{split}(T) = \frac{N_1}{N} gini(T_1) + \frac{N_2}{N} gini(T_2)$$

- For every attribute, each of the possible binary splits is considered. For a discrete-valued attribute, the attribute providing smallest $gini_{split}(T)$ is picked out to split the node. For continuous-valued properties, each possible split-point must be believed.
- The Gini impurity measures the probability of a misclassification if a label is randomly selected using the probability distribution of the branch. If all the samples belong to a single category, then index reaches its minimum (0,0).

2 Cross-entropy Impurity Index

Cross-entropy Impurity Index is based on Information theory.

- When all the samples belonging to one class are represented in a split, then null values are assumed. But if there is uniform distribution, then it is maximum.
- It permits to select the split which actually minimizes uncertainty about classification.

$$\text{Entropy} = - \sum_j p_j \log_2 p_j$$

- Information Gain (IG), determine which attribute in a given set of training feature vectors is most useful. We will use it to decide the ranking of attributes in the nodes of a decision tree.

The Information Gain (IG) can be defined as follows:

- All attributes are considered to be categorical.
- It can be adapted for continuous-valued attributes.
- The attribute which has the highest information gain is selected for split.
- Assume there are two classes, P and N.
- Consider S samples, out of these p samples belongs to class P and n samples belongs to class N.
- The amount of information, needed to decide if random example in S belongs to P or N is defined as

$$I(p, n) = -\frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$

Assume that using attribute A, a set S will be partitioned into sets $\{S_1, S_2, \dots, S_k\}$.

- If S_j contains p_j examples of P and n_j examples of N, the entropy, or the expected information needed to classify objects in all sub trees S_j is

$$E(A) = \sum_{j=1}^k \frac{p_j + n_j}{p + n} I(p_j, n_j)$$

Entropy (E)

- The "specified amount of information (in bits) needed to assign a class to a randomly drawn object in S under the optimal, shortest-length code."
- Calculate Information gain i.e. $gain(A) :$ Measures reduction in entropy, achieved because of the split. Take the split that achieves most reductions (maximizes GAIN)

$$Gain(A) = I(p, n) - E(A)$$

3 Misclassification Impurity Index

- The misclassification impurity is the simplest index, defined as:
- $I_e = 1 - \max\{p(i|i)\}$
- It is a useful criterion for pruning but not recommended for growing a decision tree since it is less sensitive to changes in the class probabilities of the nodes.

Syllabus Topic: Feature Importance

5.1.2 Feature Importance

- To predict the output value, feature importance helps to evaluate significance of every feature of multi-dimensional dataset.
- Decision trees offer a different approach based on the impurity reduction determined by every single feature.
- In particular, considering a feature x_i , its importance can be determined as:

$$\text{Importance}(x_i) = \sum_k \frac{N_k}{N} \Delta I_k$$

- The sum is extended to all nodes where x_i is used, and N_k is the number of samples reaching the node k.

Sample code for feature importance in python

```
from sklearn import tree
import DecisionTreeClassifier
from sklearn.datasets import load_iris
iris = load_iris()
clf = DecisionTreeClassifier(max_depth=4, random_state=0)
clf.fit(iris.data, iris.target)

# Output: Showing that feature 0 is name of leaf feature
# [0] ==> 0.174593813045455
# [1] ==> 0.174593813045455
# [2] ==> 0.174593813045455
# [3] ==> 0.174593813045455
```


5.7 Machine Learning (SPRU - Exam B - Comp.)

5.7

Decision Trees and Ensemble Learning

- The bagged classifier M^* uses the *voting method* i.e., the sample tuple Y is assigned the class with the most votes to tuple Y .
- For continuous values, it can be used for prediction by taking the average of all predictions for a given sample.
- Boosted tree**
 - It is a family of algorithms that are able to convert weak learners to strong learners.
 - It is a sequential ensemble method where the base learners are generated sequentially (e.g. AdaBoost).
 - The basic motivation of sequential methods is to **exploit the dependence between the base learners**. The overall performance can be boosted by weighing previously mislabelled examples with higher weight.
 - In boosting, each training tuple has weight.
 - n number of classifiers are learned iteratively.
- After learning of M_i classifier, every time the weights are updated for next classifier learning i.e. M_{i+1} . So if the tuples which were misclassified by M_i will get higher weight for next classifier.
- Use voting method, where check the votes of each classifier to get the final M^* which helps to get the accuracy.
- The extended boosting algorithm works for the prediction of continuous values.
- Boosting tends to accomplish greater accuracy as compared to bagging, there is a risk of overfitting the model.

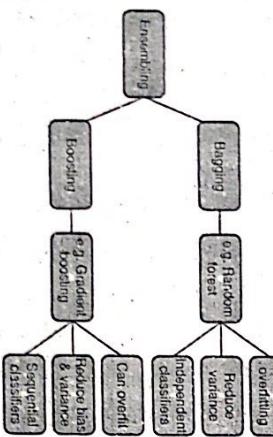


Fig. 5.1.5 : Ensemble Learning

Syllabus Topic : Random Forest

5.1.4(A) Random Forest

- It is a supervised classification algorithm.
- It creates the forest with a number of trees as the name suggests.
- If the number of trees are more, the accuracy result is high.
- Random forest handles the missing values.
- Instead of searching for the most important feature while splitting a node, it searches for the best feature among a random subset of features.

5.8 Machine Learning (SPRU - Sem B - Core)

5.8

Decision Trees and Ensemble Learning

- The algorithm selects random subset of features to split the node.
- Random Forest pseudo code**
 - Select "k" features randomly from total "m" features where $k \ll m$.
 - Using the best split point, calculate the node d from selected k features.
 - Using the best split, split the node d into child nodes.
 - Repeat 1 to 3 steps until "n" number of nodes has been reached.
 - Build forest by repeating steps 1 to 4 for "n" number times to create "n" number of trees.

Random forest prediction pseudo code

- Predict and store the outcome using the rules of each randomly created decision tree.
 - Count on the votes for each predicted target.
 - The last prediction is selected by taking the high voted predicted outcome.
- With two trees, you can see how a random forest would look like:

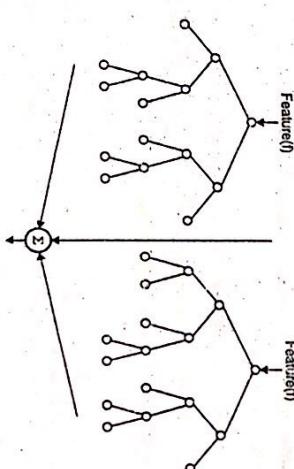


Fig. 5.1.6 : Random forest with 2 trees

Advantages

- It can be used for both regression and classification tasks.
- It is a very handy and easy to use algorithm as its default hyper parameters often produce a good prediction result.
- The classifier won't overfit the model if there are enough trees in the forest.

Disadvantages

- Many trees are generated which makes algorithm slow and not suited for real time prediction.
- It's not a descriptive tool, but only predictive model.
- If dataset is noisy, then random forest may overfit.
- For data, including categorical variables with different number of levels, random forests are biased in favour of those attributes with more levels. Therefore, the variable importance scores from random forest are not reliable for this type of data.

5.1.4(B) Ada Boost

- It combines weak classifier algorithm to form strong classifier, then by selecting, training set at every iteration multiple classifiers are combined which gives good accuracy score.
- Correct amount of weight can be assigned in final voting, which improves accuracy.

Algorithm

Step 1 : Choose the training set and train the algorithm.

Step 2 : Retrains the algorithm iteratively by selecting another set of training data based on the accuracy of previous training.

Step 3 : The weight-age depends on the accuracy achieved for each trained classifier.

Step 4 : It assigns weight to each training item.

Step 5 : Higher weights are assigned to misclassified item so that those items can appear in the training subset of next classifier with higher probability.

Step 6 : After training weights is assigned to each classifier also based on accuracy.

Step 7 : Higher weights are assigned to more accurate classifier to get more impact on the final outcome.

Step 8 : The mathematical formula for AdaBoost is given below.

Where,

$$H(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right)$$

T : Number of classifiers

$h_t(x)$: Output of weak classifier t for input x

α_t : Weight assigned to the classifier.

α_t is calculated as follows:

$\alpha_t = 0.5 * \ln((1-E)/E)$ where E is error rate.

Step 9 : Initially, all the input training examples has equal weightage. The mathematical formula to update the weight of each training example is given below.

$$D_{t+1}(i) = \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

Where,

D_t : weight at previous level.

Z_t : sum of all weights used to normalize the weights.

y_i is y coordinate of training example (x_i, y_i) for simplicity.

5.1.4(C) Gradient Tree Boosting

- Gradient tree boosting is a technique that allows you to build a tree ensemble step by step with the goal of minimizing a target loss function. Gradient boosting has main three components:

1. A loss function to be optimized.
2. A weak learner to make predictions.
3. An additive model to add weak learners to minimize the loss function.

Mean squared error (MSE) as loss defined as:

$$\text{Loss} = \text{MSE} = \sum (y_i - \hat{y}_i)^2$$

Where, y_i = i^{th} target value, \hat{y}_i = i^{th} prediction, $L(y, \hat{y})$ is loss function

MSE should be minimum for our predictions.

Use a gradient descent and update our predictions based on a learning rate, we can find the values where MSE is minimized.

$$\hat{y}_i = y_i + \alpha * \delta \sum (y_i - \hat{y}_i)^2 / \delta y_i$$

Which becomes,

$$\hat{y}_i = y_i - \alpha * 2 * \sum (y_i - \hat{y}_i)$$

Where, α is learning rate and $\sum (y_i - \hat{y}_i)$ is the sum of residuals.

5.1.4(D) Voting Classifier

Step 1 : Assume mean is the prediction of all variables.

Step 2 : For each observation calculate errors using mean (latest prediction).

Step 3 : Find the variable that can split the errors perfectly and find the value for the split. This is assumed to be the latest prediction.

Step 4 : Calculate errors of each observation from the mean of both the sides of the split (latest prediction).

Step 5 : Repeat the step 3 and 4 till the objective function maximizes/minimizes.

Step 6 : Take a weighted mean of all the classifiers to come up with the final model.

Syllabus Topic : Voting Classifier

5.1.4(D) Voting Classifier

Voting classifier has two strategies

1. Hard Voting : Final class label is predicted based on the most frequent class label of classification models.
2. Soft Voting : Final class label is predicted by averaging the class probabilities.

y_i is y coordinate of training example (x_i, y_i) for simplicity.

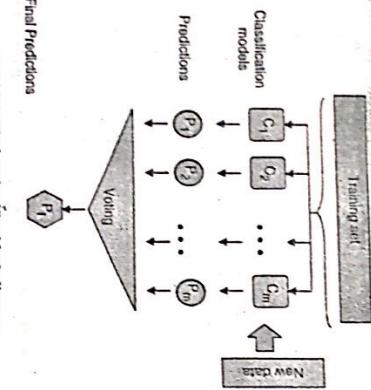


Fig 5.1.7: Voting classifier block diagram

or Hard Voting or Majority VotingPredict the class label Y via majority voting of each classifier C_j :

$$Y = \text{mode} \{C_1(x), C_2(x), \dots, C_m(x)\}$$

Assume 3 classifiers and predicted class of each as given below:

classifier 1 > class 0

classifier 2 > class 0

classifier 3 > class 1

$$y = \text{mode} \{0, 0, 1\} = 0$$

So final predicted class for the sample is "0"

or Weighted Majority VoteWeighted majority vote can be calculated by associating a weight w_j with classifier C_j :

$$\hat{y} = \arg \max_i \sum_{j=1}^m w_j X_\lambda(C_j(x) = i)$$

where X_λ is the characteristic

A is the set of unique class labels.

- classifier 1 > class 0
- classifier 2 > class 0
- classifier 3 > class 1

Assign the weights as {0.2, 0.2, 0.6} respectively for each class.

$$\text{Prediction for } Y_{15} \\ \arg \max_j [0.2 \times i_6 + 0.2 \times i_9 + 0.6 \times i_1] = 1$$

or Soft VotingClass label is predicted based on the predicted probabilities p for the classifier. This works only if classifiers are well-calibrated.

$$\hat{y} = \arg \max_i \sum_{j=1}^m w_j p_{ij}$$

where W_j is the weight that can be assigned to the j th classifier.

Assume the same example with binary classification having class labels {0, 1} our ensemble could make the following prediction:

$$C1(x) \rightarrow [0.9, 0.1]$$

$$C2(x) \rightarrow [0.8, 0.2]$$

$$C3(x) \rightarrow [0.4, 0.6]$$

Using uniform weights, we compute the average probabilities:

$$p(i_6 | x) = \frac{0.9 + 0.8 + 0.4}{3} = 0.7$$

$$p(i_1 | x) = \frac{0.1 + 0.2 + 0.6}{3} = 0.3$$

$$\hat{y} = \arg \max_i [p(i_6 | x), p(i_1 | x)] = 0$$

However, assigning the weights {0.1, 0.1, 0.8} would yield a prediction $\hat{y} = 1$:

$$p(i_6 | x) = 0.1 \times 0.9 + 0.1 \times 0.8 + 0.8 \times 0.4 = 0.49$$

$$p(i_1 | x) = 0.1 \times 0.1 + 0.2 \times 0.1 + 0.8 \times 0.6 = 0.51$$

$$\hat{y} = \arg \max_i [p(i_6 | x), p(i_1 | x)] = 1$$

Syllabus Topic : Clustering Fundamentals - Basics

5.2 Clustering Fundamentals

5.2.1 Basics of Clustering

- Clustering is an unsupervised learning problem.
- Data is divided into different groups, called as clusters.
- No classes or groups are known earlier.
- Clusters should be made based on data objects having high inter similarity and low intra similarity.
- In hard clustering techniques, each element must belong to a single cluster and in soft clustering or fuzzy clustering, is based on a membership score that defines how much the elements are "compatible" with each cluster.

- Graphical representation can be shown:

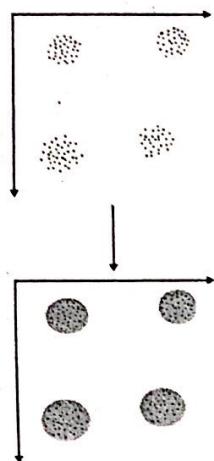


Fig 5.2.1 : Clustering Graphical example

- From Fig. 5.2.1, we can observe there are 4 clusters based on their geometrical distance grouped together is called distance based clustering.

- Conceptual clustering is based on common or descriptive concepts of data objects.

Applications

Clustering algorithms can be applied in many disciplines:

- **Marketing :** Clustering can be used for targeted marketing. For e.g. Given a customer database, containing properties and past buying records. Similar groups of customers can be identified and grouped into one cluster.

- **Biology :** Clustering can also be employed in classifying plants and animals into different categories based on their characteristics.

Applications

Insurance : Different groups of policy holders can be placed with clustering. For e.g. Identify policy fraud or customers with high claim.

- City planning : Groups of houses based on location and type.

- Earthquake studies : Clustering can also be used to identify dangerous zones based on earthquake epicentre.

- WWW : Clustering can be used to find groups of similar access patterns using weblog data. It can also be used for classification of documents.

Syllabus Topic : K-Means

5.2.2 K-means

- In 1967, J. MacQueen and then in 1973 J. A. Hartigan and M. A. Wong developed K-means clustering algorithm.

- In k-means, k is the number of clusters given by user and objects are classified into k clusters based on their attributes.

- K-means is one of the simplest unsupervised learning algorithms.

- Define K centroids for K clusters which are generally far away from each other.

- Group the objects into clusters based on the distance with respect to centroid

- After this first step, again calculate the new centroid for each cluster based on the elements of that cluster.
- Follow the same method and group the elements based on new centroid.
- At every step, the centroid changes and elements move from one cluster to another.
- Do the same process till no element is moving from one cluster to another i.e. till two consecutive steps with the same centroid and the same elements are obtained.
- Finally, this algorithm aims at minimizing an objective function, in this case a squared error function. The objective function is given below,

$$J = \sum_{j=1}^k \sum_{i=1}^n \|x_i^{(j)} - c_j\|^2$$

Where, $x_i^{(j)}$ = A data point

c_j = The cluster centre

n = Number of data points

k = Number of clusters

$$\|x_i^{(j)} - c_j\|^2 = \text{Distance measure between a data point } x_i^{(j)} \text{ and the cluster centre } c_j$$

5.2.2(A) K-means Algorithm

k: number of clusters

x_1, x_2, \dots, x_n : sample feature vectors x_1, x_2, \dots, x_n

m_i : the mean of the vectors in cluster i

Assume $k < n$.

Make initial guesses for the means m_1, m_2, \dots, m_k .

Until there are no changes in any mean.

Use the estimated means to classify the samples into clusters.

for $i = 1$ to k

Replace m_i with the mean of all of the samples for cluster i

end_for

end_until

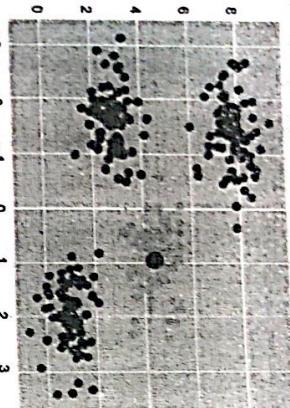
Following the three steps are repeated until convergence:

Iterate till no object moves to a different group :

Step 1 : Find the centroid coordinate.

Step 2 : Find the distance of each object to the centroids.

Step 3 : Based on minimum distance group the objects.



Strength of K-means clustering

Relatively efficient : $O(kn)$.

where n is number of objects,

k is number of clusters,

l is number of iterations.

Normally, $k < n$.

K-means often terminates at a local optimum.

Techniques like deterministic annealing and genetic algorithms are used to get the global optimum solution.

Weakness of K-means clustering

Applicable only when means are defined.

Need to specify k , the number of clusters, in advance.

Unable to handle noisy data and outliers (outlier: objects with extremely large values)

Not suitable to discover clusters with non-convex shapes.

Syllabus Topic : Finding Optimal Number of Clusters

6.2.2(C) Finding Optimal Number of Clusters

The disadvantage of k-means is to find the optimal number of clusters. If the number of clusters are less , then large element grouping can be formed with heterogeneous elements. With the more number of clusters , it will be difficult to get the dissimilarities among clusters.

Therefore, different methods to determine the number of clusters are

1. Optimizing the inertia
2. Silhouette score
3. Calinski-Harabasz Index
4. Cluster instability

1. Optimizing the Inertia

- The assumption that an appropriate number of clusters must produce a small inertia.
- When number of clusters is same as number of data elements, this value reaches its minimum (0.0).
- So don't look for the minimum, but for a value which is a trade-off between the inertia and the number of clusters.
- This method calculates the inertia for different number of clusters.

2. Silhouette score

Silhouette score is a way to measure how the elements within a cluster are close to the points in its neighbouring clusters.

- It is based on the principle of "maximum internal cohesion and maximum cluster separation" means how similarly an object is in its own cluster (cohesion) compared to other clusters (separation).
- It finds the optimal value of k (number of clusters) during clustering.

Define a distance metric and compute the average intra-cluster distance for each element :

$$a(\bar{x}_i) = \sum_{j \in C} d(\bar{x}_i, \bar{x}_j) \quad \forall \bar{x}_i \in C$$

Calculate the average nearest-cluster distance which is the lowest inter-cluster distance.

$$b(\bar{x}_i) = \min_{j \in D \setminus C} d(\bar{x}_i, \bar{x}_j) \quad \forall \bar{x}_i \in C \text{ where } D = \arg\min_i d(C, D)$$

The silhouette score for an element x_i is defined as :

$$S(\bar{x}_i) = \frac{b(\bar{x}_i) - a(\bar{x}_i)}{\max(a(\bar{x}_i), b(\bar{x}_i))}$$

This value of silhouette score is bounded between -1 and 1.

Score closer to 1 means assigned to the cluster correctly and score closer to -1 is assigned to a wrong cluster. A score close to 0 means the point lies between almost at the boundary of both the clusters.

scikit-learn allows computing the average silhouette score to have an immediate overview for different numbers of clusters : The program is taken from

https://scikit-learn.org/stable/auto_examples/cluster/plot_kmeans_silhouette_analysis.html

```
from sklearn import datasets
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.cluster import KMeans
from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler
```

The particular writing is over different clusters and 3 clusters placed close

```
# k-means
X = make_blobs(n_samples=50)

n_features = 2
centers = 3
cluster_std = 1
center_box = [10, 50, 100]

# n_clusters = 3
random_state = 1 # For reproducibility

# For 3 clusters in range n_clusters
# Create subplot with 1 row and 2 columns
f, (ax1, ax2) = plt.subplots(1, 2)

# For 1st subplot, the silhouette coefficient can range from -1.0 to 1.0. In this example it's all
# 1s within [-0.1, 1]
ax1.set_xlim([-0.1, 1])

# The (n_clusters-1)*10 is for inserting blank square between silhouette plots
# of individual clusters to delineate them clearly.
ax1.set_ylim([0, len(X) + (n_clusters + 1) * 10])

# Initialize the cluster with n_clusters value and a random generator
# seed of 10 for reproducibility
clusters = KMeans(n_clusters=n_clusters, n_init=10,
                  random_state=10)

# Plotting silhouette scores for all the samples
# The silhouette score for a sample is computed as the difference between
# a sample's average distance to all the other samples and the second
# closest cluster's average distance to all the other samples, divided by the
# maximum of those two values. The silhouette scores therefore lie in the
# range [-1, 1]. They are higher for more similar samples in a cluster, lower
# for rather dissimilar samples belonging to the same cluster, and intermediate
# values for intermediate cases.
# Compute the silhouette scores for all the samples
# silhouette = sample_silhouette_values[cluster_labels == i]
# i-th cluster silhouette values
# facecolor = color if cluster_label == i else "black"
# color = "#%02x%02x%02x" % (np.random.randint(0, 256), np.random.randint(0, 256),
#                             np.random.randint(0, 256))

# And plot showing the actual clusters formed
colors = cm.cmap_color_cycle(cluster_labels.astype(float)/n_clusters)
ax2.scatter(X[:, 0], X[:, 1], marker='o', s=30, c=colors, alpha=0.7,
            edgecolor="black")

# Plotting the clusters
centers = cluster_centers_.copy()
# Draw white circles at cluster centers
ax2.scatter(centers[:, 0], centers[:, 1], marker='^', s=100, c='white',
            zorder=10, edgecolor="black")

# Drawing horizontal ellipses at cluster centers
for i, center in enumerate(centers):
    ax2.ellipsize(xy=center, width=10, height=10, angle=0)
```

For 3 clusters in range n_clusters

Create subplot with 1 row and 2 columns

f, (ax1, ax2) = plt.subplots(1, 2)

For 1st subplot, the silhouette coefficient can range from -1.0 to 1.0. In this example it's all

1s within [-0.1, 1]

ax1.set_xlim([-0.1, 1])

The (n_clusters-1)*10 is for inserting blank square between silhouette

plots of individual clusters to delineate them clearly.

ax1.set_ylim([0, len(X) + (n_clusters + 1) * 10])

Initialize the cluster with n_clusters value and a random generator

seed of 10 for reproducibility

clusters = KMeans(n_clusters=n_clusters, n_init=10,

random_state=10)

The silhouette score for a sample is computed as the difference between

a sample's average distance to all the other samples and the second

closest cluster's average distance to all the other samples, divided by the

maximum of those two values. The silhouette scores therefore lie in the

range [-1, 1]. They are higher for more similar samples in a cluster, lower

for rather dissimilar samples belonging to the same cluster, and intermediate

values for intermediate cases.

Compute the silhouette scores for all the samples

silhouette = sample_silhouette_values[cluster_labels == i]

i-th cluster silhouette values

facecolor = color if cluster_label == i else "black"

color = "#%02x%02x%02x" % (np.random.randint(0, 256), np.random.randint(0, 256),

np.random.randint(0, 256))

And plot showing the actual clusters formed

colors = cm.cmap_color_cycle(cluster_labels.astype(float)/n_clusters)

ax2.scatter(X[:, 0], X[:, 1], marker='o', s=30, c=colors, alpha=0.7,

edgecolor="black")

Plotting horizontal ellipses at cluster centers

for i, center in enumerate(centers):

ax2.ellipsize(xy=center, width=10, height=10, angle=0)

4 clusters and sort them.

ith cluster silhouette values = sample_silhouette_values[cluster_labels == i]

size cluster[j] = (ith cluster silhouette values).sum()

gamma = 1 - lower + size cluster[i]

color = "#%02x%02x%02x" % (np.random.randint(0, 256), np.random.randint(0, 256),

np.random.randint(0, 256))

ax1.fill_betweenx(gamma, 0, cluster_silhouette_values.sum(), color)

ith cluster silhouette values

facecolor = color if cluster_label == i else "black"

color = "#%02x%02x%02x" % (np.random.randint(0, 256), np.random.randint(0, 256),

np.random.randint(0, 256))

And plot showing the actual clusters formed

colors = cm.cmap_color_cycle(cluster_labels.astype(float)/n_clusters)

ax2.scatter(X[:, 0], X[:, 1], marker='o', s=30, c=colors, alpha=0.7,

edgecolor="black")

Plotting the clusters

centers = cluster_centers_.copy()

Draw white circles at cluster centers

ax2.scatter(centers[:, 0], centers[:, 1], marker='^', s=100, c='white',

zorder=10, edgecolor="black")

Axes object

1 asking the clusters

centers = cluster_centers_.copy()

Draw white circles at cluster centers

ax2.scatter(centers[:, 0], centers[:, 1], marker='^', s=100, c='white',

zorder=10, edgecolor="black")

Plot showing the actual clusters formed

colors = cm.cmap_color_cycle(cluster_labels.astype(float)/n_clusters)

ax2.scatter(X[:, 0], X[:, 1], marker='o', s=30, c=colors, alpha=0.7,

edgecolor="black")

Plotting the clusters

centers = cluster_centers_.copy()

Draw white circles at cluster centers

ax2.scatter(centers[:, 0], centers[:, 1], marker='^', s=100, c='white',

zorder=10, edgecolor="black")

Axes object

1 asking the clusters

centers = cluster_centers_.copy()

Draw white circles at cluster centers

ax2.scatter(centers[:, 0], centers[:, 1], marker='^', s=100, c='white',

zorder=10, edgecolor="black")

Plot showing the actual clusters formed

colors = cm.cmap_color_cycle(cluster_labels.astype(float)/n_clusters)

ax2.scatter(X[:, 0], X[:, 1], marker='o', s=30, c=colors, alpha=0.7,

edgecolor="black")

Plotting the clusters

centers = cluster_centers_.copy()

Draw white circles at cluster centers

ax2.scatter(centers[:, 0], centers[:, 1], marker='^', s=100, c='white',

zorder=10, edgecolor="black")

Axes object

1 asking the clusters

centers = cluster_centers_.copy()

Draw white circles at cluster centers

ax2.scatter(centers[:, 0], centers[:, 1], marker='^', s=100, c='white',

zorder=10, edgecolor="black")

Plot showing the actual clusters formed

colors = cm.cmap_color_cycle(cluster_labels.astype(float)/n_clusters)

ax2.scatter(X[:, 0], X[:, 1], marker='o', s=30, c=colors, alpha=0.7,

edgecolor="black")

Plotting the clusters

centers = cluster_centers_.copy()

Draw white circles at cluster centers

ax2.scatter(centers[:, 0], centers[:, 1], marker='^', s=100, c='white',

zorder=10, edgecolor="black")

Plot showing the actual clusters formed

colors = cm.cmap_color_cycle(cluster_labels.astype(float)/n_clusters)

ax2.scatter(X[:, 0], X[:, 1], marker='o', s=30, c=colors, alpha=0.7,

edgecolor="black")

Plotting the clusters

centers = cluster_centers_.copy()

Draw white circles at cluster centers

ax2.scatter(centers[:, 0], centers[:, 1], marker='^', s=100, c='white',

zorder=10, edgecolor="black")

Plot showing the actual clusters formed

colors = cm.cmap_color_cycle(cluster_labels.astype(float)/n_clusters)

ax2.scatter(X[:, 0], X[:, 1], marker='o', s=30, c=colors, alpha=0.7,

edgecolor="black")

Plotting the clusters

centers = cluster_centers_.copy()

Draw white circles at cluster centers

ax2.scatter(centers[:, 0], centers[:, 1], marker='^', s=100, c='white',

zorder=10, edgecolor="black")

Plot showing the actual clusters formed

colors = cm.cmap_color_cycle(cluster_labels.astype(float)/n_clusters)

ax2.scatter(X[:, 0], X[:, 1], marker='o', s=30, c=colors, alpha=0.7,

edgecolor="black")

Plotting the clusters

centers = cluster_centers_.copy()

Draw white circles at cluster centers

ax2.scatter(centers[:, 0], centers[:, 1], marker='^', s=100, c='white',

zorder=10, edgecolor="black")

Plot showing the actual clusters formed

colors = cm.cmap_color_cycle(cluster_labels.astype(float)/n_clusters)

ax2.scatter(X[:, 0], X[:, 1], marker='o', s=30, c=colors, alpha=0.7,

edgecolor="black")

Plotting the clusters

centers = cluster_centers_.copy()

Draw white circles at cluster centers

ax2.scatter(centers[:, 0], centers[:, 1], marker='^', s=100, c='white',

zorder=10, edgecolor="black")

Plot showing the actual clusters formed

colors = cm.cmap_color_cycle(cluster_labels.astype(float)/n_clusters)

ax2.scatter(X[:, 0], X[:, 1], marker='o', s=30, c=colors, alpha=0.7,

edgecolor="black")

Plotting the clusters

centers = cluster_centers_.copy()

Draw white circles at cluster centers

ax2.scatter(centers[:, 0], centers[:, 1], marker='^', s=100, c='white',

zorder=10, edgecolor="black")

Plot showing the actual clusters formed

colors = cm.cmap_color_cycle(cluster_labels.astype(float)/n_clusters)

ax2.scatter(X[:, 0], X[:, 1], marker='o', s=30, c=colors, alpha=0.7,

edgecolor="black")

Plotting the clusters

centers = cluster_centers_.copy()

Draw white circles at cluster centers

ax2.scatter(centers[:, 0], centers[:, 1], marker='^', s=100, c='white',

zorder=10, edgecolor="black")

Plot showing the actual clusters formed

colors = cm.cmap_color_cycle(cluster_labels.astype(float)/n_clusters)

ax2.scatter(X[:, 0], X[:, 1], marker='o', s=30, c=colors, alpha=0.7,

edgecolor="black")

Plotting the clusters

centers = cluster_centers_.copy()

Draw white circles at cluster centers

ax2.scatter(centers[:, 0], centers[:, 1], marker='^', s=100, c='white',

zorder=10, edgecolor="black")

Plot showing the actual clusters formed

colors = cm.cmap_color_cycle(cluster_labels.astype(float)/n_clusters)

ax2.scatter(X[:, 0], X[:, 1], marker='o', s=30, c=colors, alpha=0.7,

edgecolor="black")

Plotting the clusters

centers = cluster_centers_.copy()

Draw white circles at cluster centers

ax2.scatter(centers[:, 0], centers[:, 1], marker='^', s=100, c='white',

zorder=10, edgecolor="black")

Plot showing the actual clusters formed

colors = cm.cmap_color_cycle(cluster_labels.astype(float)/n_clusters)

ax2.scatter(X[:, 0], X[:, 1], marker='o', s=30, c=colors, alpha=0.7,

edgecolor="black")

Plotting the clusters

centers = cluster_centers_.copy()

Draw white circles at cluster centers

ax2.scatter(centers[:, 0], centers[:, 1], marker='^', s=100, c='white',

zorder=10, edgecolor="black")

Plot showing the actual clusters formed

colors = cm.cmap_color_cycle(cluster_labels.astype(float)/n_clusters)

ax2.scatter(X[:, 0], X[:, 1], marker='o', s=30, c=colors, alpha=0.7,

edgecolor="black")

Plotting the clusters

centers = cluster_centers_.copy()

Draw white circles at cluster centers

ax2.scatter(centers[:, 0], centers[:, 1], marker='^', s=100, c='white',

zorder=10, edgecolor="black")

Plot showing the actual clusters formed

colors = cm.cmap_color_cycle(cluster_labels.astype(float)/n_clusters)

ax2.scatter(X[:, 0], X[:, 1], marker='o', s=30, c=colors, alpha=0.7,

edgecolor="black")

Plotting the clusters

centers = cluster_centers_.copy()

Draw white circles at cluster centers

ax2.scatter(centers[:, 0], centers[:, 1], marker='^', s=100, c='white',

zorder=10, edgecolor="black")

Plot showing the actual clusters formed

colors = cm.cmap_color_cycle(cluster_labels.astype(float)/n_clusters)

ax2.scatter(X[:, 0], X[:, 1], marker='o', s=30, c=colors, alpha=0.7,

edgecolor="black")

Plotting the clusters

centers = cluster_centers_.copy()

Draw white circles at cluster centers

ax2.scatter(centers[:, 0], centers[:, 1], marker='^', s=100, c='white',

zorder=10, edgecolor="black")

Plot showing the actual clusters formed

colors = cm.cmap_color_cycle(cluster_labels.astype(float)/n_clusters)

ax2.scatter(X[:, 0], X[:, 1], marker='o', s=30, c=colors, alpha=0.7,

edgecolor="black")

Plotting the clusters

centers = cluster_centers_.copy()

Draw white circles at cluster centers

ax2.scatter(centers[:, 0], centers[:, 1], marker='^', s=100, c='white',

zorder=10, edgecolor="black")

Plot showing the actual clusters formed

colors = cm.cmap_color_cycle(cluster_labels.astype(float)/n_clusters)

ax2.scatter(X[:, 0], X[:, 1], marker='o', s=30, c=colors, alpha=0.7,

edgecolor="black")

Plotting the clusters

17. K-means Clustering (Spiral, Part 3 - Cont.)	8/23	Practical Tools and Examples (Assuming)
<ul style="list-style-type: none"> - Consider k = 2 clusters with $\mu_1(7, 1)$, $\mu_2(17, 1)$. Between two clusterings with the same number (k) of clusters, the probability is highest for having distance between samples of clusterings of 1. (or minimum). 	$\mu_1 = (7, 1), \mu_2 = (17, 1), C_1, C_2, M$	

- We need to find the value of k that minimizes $\|C\|_F^2$ (sum of squares for C).

Spiral Data: DBSCAN

DBSCAN: Density Based Spatial Clustering of Applications with Noise

The algorithm DBSCAN works on the formal notion of density reachability for k points and points are designated by ϵ -neighbor clusters of arbitrary shape. The outcome of the algorithm is of the form C_1, C_2, \dots, C_n where clusters are efficiently supported by spatial index structures, i.e. at least in moderately dimensional spaces.

Explanation of DBSCAN Steps

- Equation (17.5) and Minimum points (17.6)(a) are the two parameters selected by DBSCAN. $\epsilon > 0$ is a fixed point as chosen. δ is the radius of the neighborhood of a point. If all the points within δ are neighbors, then they are called ϵ -neighbors.
- The distance between the starting point and its neighbors is calculated. If the points within δ are all ϵ -neighbors, then the starting point is marked as visited. The process continues until no new point is found.
- The above steps are then repeated for all the remaining neighborhoods.
- If all the points within each of the neighborhoods are less than δ , then they are marked as noise.
- If all the points within each of the clusters are visited then the algorithm proceeds by selecting other ϵ -neighborhoods.
- For any cluster, we have:
 - A central point (p) is core point.
 - A distance from the core point (δ_{cp}).
 - Minimum number of points within the specified distance (M_{cp}).
- **Properties**
 - One of the limitations is that, for very low to very high M_{cp} and δ the clustering may become a poor representation of the dataset. Sometimes at high dimensional data set, it is difficult to decide.
 - Some of the cluster distributions may be highly nonconvex. If a portion of the area in the cluster may be very sparse compared to the other area, some of the neighbors may not have clusters or noise may be present.
 - The process is extremely sensitive to noise which leads to very different clusters.
- **Code for DBSCAN**: Implementing DBSCAN and understanding stampyplot (https://github.com/stampylion/StampyPlot)
- **Advantages**
 - (i) Discover clusters of arbitrary shape.
 - (ii) Handles noise.
 - (iii) One pass.
 - (iv) Short memory requirements as termination condition.

18. Hierarchical Clustering (Spiral, Part 4 - Cont.)	8/24	Practical Tools and Examples (Assuming)
<ul style="list-style-type: none"> - Clusters are formed in a hierarchical manner. - Clusters of arbitrary shape and size are formed which are dense. - The algorithm is as follows: <ol style="list-style-type: none"> (i) Core objects and density reachable objects are found out, where these density reachable core objects and their clusters are discovered. (ii) When no new points can be added to any of the clusters then execution may be stopped. (iii) Clusters are formed by regions of low density (noise). - Clusters will not affect the creation of clusters. 	DBSCAN Method Step 1: Select ϵ and δ Step 2: Find core points Step 3: Assign each point as the cluster from which the nearest point is at least δ distance away. Step 4: Assign non-core points within ϵ distance of a point. Step 5: Assign non-core points within δ distance of a point. Step 6: Assign noise.	

- Clusters are discovered.
- When no new points can be added to any of the clusters then execution may be stopped.
- Clusters are formed by regions of low density (noise).
- Clusters will not affect the creation of clusters.

Spiral Data: DBSCAN

DBSCAN: Density Based Spatial Clustering of Applications with Noise

The algorithm DBSCAN works on the formal notion of density reachability for k points and points are designated by ϵ -neighbor clusters of arbitrary shape. The outcome of the algorithm is of the form C_1, C_2, \dots, C_n where clusters are efficiently supported by spatial index structures, i.e. at least in moderately dimensional spaces.

Explanation of DBSCAN Steps

- Equation (17.5) and Minimum points (17.6)(a) are the two parameters selected by DBSCAN. $\epsilon > 0$ is a fixed point as chosen. δ is the radius of the neighborhood of a point. If all the points within δ are neighbors, then the starting point is marked as visited. The process continues until no new point is found.
- The above steps are then repeated for all the remaining neighborhoods.
- If all the points within each of the neighborhoods are less than δ , then they are marked as noise.
- If all the points within each of the clusters are visited then the algorithm proceeds by selecting other ϵ -neighborhoods.
- For any cluster, we have:
 - A central point (p) is core point.
 - A distance from the core point (δ_{cp}).
 - Minimum number of points within the specified distance (M_{cp}).
- **Properties**
 - One of the limitations is that, for very low to very high M_{cp} and δ the clustering may become a poor representation of the dataset. Sometimes at high dimensional data set, it is difficult to decide.
 - Some of the cluster distributions may be highly nonconvex. If a portion of the area in the cluster may be very sparse compared to the other area, some of the neighbors may not have clusters or noise may be present.
 - The process is extremely sensitive to noise which leads to very different clusters.
- **Code for DBSCAN**: Implementing DBSCAN and understanding stampyplot (https://github.com/stampylion/StampyPlot)
- **Advantages**
 - (i) Discover clusters of arbitrary shape.
 - (ii) Handles noise.
 - (iii) One pass.
 - (iv) Short memory requirements as termination condition.

Machine Learning (GPU - Sem B - Comp.)

5.25

Decision Trees and Ensemble Learning

```

X = np.random.rand(100, 2)
y = np.sin(6.28 * X[0] + 0.1 * np.random.randn(100))
# Compute DSCAN
dscn = DSCAN(X, y, 0.3, max_iter=10, init='k-means++', sample_size=10,
              random_state=42, transform=None)

# Number of clusters in X, ignoring noise if present.
n_clusters = len(dscn.labels_)

# Adjusted Rand Index
adjusted_rand_index_score(dscn.labels_, y)

# Estimated number of clusters in X, ignoring noise if present.
estimated_n_clusters = len(np.unique(dscn.labels_))

# Estimated number of noise points
len(np.where(dscn.labels_ == -1)[0])

# Estimated number of true points
len(np.where(dscn.labels_ >= 0)[0])

```

Machine Learning (GPU - Sem B - Comp.)

5.26

Decision Trees and Ensemble Learning

```

# K-means clustering example
from sklearn.cluster import KMeans
kmeans = KMeans(n_clusters=6, random_state=0).fit(X)

# Estimated number of clusters
estimated_n_clusters = len(np.unique(kmeans.labels_))

# Estimated number of clusters
estimated_n_clusters = len(np.unique(kmeans.labels_))

# Estimated number of noise points
len(np.where(kmeans.labels_ == -1)[0])

# Estimated number of true points
len(np.where(kmeans.labels_ >= 0)[0])

```

Output

Automatically created model for GPU in interactive environment.

Estimated number of clusters: 3

Estimated number of noise points: 18

Adjusted Rand Index: 0.952

Adjusted Rand Index (K-Means): 0.952

Estimated number of clusters: 3

Estimated number of noise points: 18

Adjusted Rand Index: 0.952

Adjusted Rand Index (K-Means): 0.952

Estimated number of clusters: 3

Syllabus topic : Spectral Clustering

5.24 Spectral Clustering

The goal of spectral clustering is to cluster data that is connected but not necessarily compact or clustered within convex boundaries.

Three main steps of spectral clustering :

1. Create a similarity graph between our N objects to cluster.
2. Compute the first k eigenvectors of its Laplacian matrix to define a feature vector for each object.
3. Run k-means on these features to separate objects into k classes.

Y values are only for first 200 samples. To get all values, change n_samples=200, cluster_id=0, random_state=42, transform=None

ValueError: Input contains NaN, infinity or a value too large for dtype('float64').

ValueError: Input contains NaN, infinity or a value too large for dtype('float64').

- Compactness, e.g. k-means, mixture models
 - Connectivity, e.g. spectral clustering

Charles H Martin (<https://calculatedcontent.com/2012/10/09/spectral-clustering/>) has described the algorithm as given

1. - project your data into R^k
 2. - define an *Affinity* matrix A , using a Gaussian Kernel K or say just an *Adjacency* matrix (i.e. $A_{ij} = \delta_{ij}$).
 3. - construct the Graph Laplacian from A . (i.e. decide on a normalization).
 4. - solve an Eigen value problem, such as $L_\tau = X_\tau$ (or a Generalized Eigen value problem $L_\tau = \lambda D_\tau$).
 5. - select k eigenvectors $\{v_i, i = 1, k\}$ corresponding to the k lowest (or highest) eigenvalues $\{\lambda_i, i = 1, k\}$, to define a k -dimensional subspace $L^T P$.

5.3 Evaluation Methods Based on Ground Truth-Homogeneity, Completeness, Adjusted Rand Index

Before the prediction of any new sample, evaluation of model is necessary. The model can be evaluated with Homogeneity, Completeness, Adjusted Rand Index.

- The homogeneity score is used to determine the optimal number of clusters. A clustering result satisfies homogeneity if all of its clusters contain only data points which are members of a single class.

Where,

C. Classification

$H(C_k)$: measure of uncertainty after clustering in determining the right class. It's necessary to normalize this value considering the initial entropy of the class set $H(C)$.

5.3 Evaluation Methods Based on Ground Truth- Homogeneity, Completeness, Adjusted Rand Index

Before the prediction of any new sample, evaluation of model is necessary. The model can be evaluated with

- ## 1. Homogeneity (h)

The homogeneity score was calculated as the percentage of all of its clusters contain only data points which are members of a single class.

$$h = 1 - \frac{H(C)}{H(C|K)}$$

C: class set

$H(C|K)$: measure of uncertainty after clustering in determining the right class, it's necessary to normalize this value considering the initial entropy of the class set $H(C)$.

considering the initial entropy of the class set H

considering the mutual entropy of the classes and

Scanned with CamScanner

Machine Learning (SPPU - Sem 8 - Comp) 5-29

Decision Trees and Ensemble Learning

```

print("Estimated number of clusters: %d (%d clusters)" % (k, n_clusters))
print("Homogeneity: %0.3f (%d genuine homogenity score[%d, true_label])")
print("Completeness: %0.3f (%d metrics completeness score[%d, true_label])")
print("V-measure: %0.3f (%d genuine v-measure score[%d, true_label])")
print("Adjusted Rand Index: %0.3f (%d adjusted rand index score[%d, true_label])")
print("Adjusted Mutual Information: %0.3f (%d adjusted mutual info score[%d, true_label])")
print("Silhouette Coefficient: %0.3f (%d silhouette score[%d, true_label])")
print("Adjusted silhouette score (%d, true_label, true_labels))")

# Compute silhouette scores for each sample, and find their average
# silhouette_score = np.mean(silhouette_samples)

# Display the results
plt.title("Silhouette analysis on a 2D dataset")
plt.xlabel("The 1st principal component")
plt.ylabel("The 2nd principal component")
plt.scatter(X[:, 0], X[:, 1], c=y)
for i in range(n_clusters):
    cluster_center = X[y == i].mean(axis=0)
    plt.plot([cluster_center[0], cluster_center[1]], [cluster_center[1], cluster_center[0]], color='black')
    plt.annotate(i, (cluster_center[0], cluster_center[1]), color='black')
plt.show()

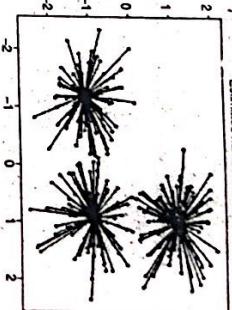
```

Output

A automatically created module for Python interactive environment
 Estimated number of clusters: 3
 Homogeneity: 0.912
 Completeness: 0.912
 V-measure: 0.912
 Adjusted Mutual Information: 0.911
 Silhouette Coefficient: 0.753

Machine Learning (SPPU - Sem 8 - Comp) 5-30

Decision Trees and Ensemble Learning



Syllabus Topic : Introduction to Meta Classifier

5.4 Introduction to Meta Classifier

Meta classifier is a classifier, which is usually a proxy to the main classifier, used to provide additional data processing. Meta classifier doesn't implement a classification algorithm on its own, but uses another classifier to do the actual work. In addition, the meta-classifier adds another processing step that is performed before the actual base-classifiers work.

Use meta classifier before the classification on both training and testing data. The Meta classifier technique is used to implement arbitrary pre-processing within cross-validation analysis.

Syllabus Topic : Concepts of Weak and Eager Learner

5.4.1 Concepts of Weak and Eager Learner

- A weak learner: It is defined to be a classifier that is only slightly correlated with the true classification. In contrast, a strong learner is a classifier that is arbitrarily well-correlated with the true classification.
- Eager learner : Given a set of training tuples, constructs a classification model before receiving new (e.g., test) data to classify.

Review Questions

- Q. 1 Explain different Impurity measures.
- Q. 2 What is Ensemble Learning?
- Q. 3 Explain following algorithms in detail
 - (a) Random Forest
 - (b) Ada Boost
 - (c) K-Means
 - (d) DBSCAN
- Q. 4 What are different methods to determine number of clusters? Explain.

Clustering Techniques

Divisive Hierarchical Clustering

In this data objects are grouped in a top down manner. Initially, all objects are in one cluster. Then the cluster is subdivided into smaller and smaller pieces, until each object forms a cluster on its own or until it satisfies certain termination conditions as the desired number of clusters are obtained. Divisive methods are not generally available, and rarely have been applied.

Difference between Agglomerative and Divisive Hierarchical Clustering

Agglomerative	Divisive
Start by considering each individual element as a cluster.	Start with only one cluster by combining all elements.
At every step, clusters are merged based on closest pair until a single or k clusters are left	At every step, the cluster is divided or split until each cluster contains a single element or k clusters are left.

Advantages

- This algorithm is simple and gives an output as a hierarchy
- The structure obtained is easy to understand and more informative.
- There is no need to pre-specify the number of clusters

Disadvantages

- Merging and splitting is critical once the clusters are formed, as undo is not possible. Every time it performs on newly generated clusters.
- If decision to merge and split is wrong, then low quality clusters are formed
- Scikit-learn implements only the agglomerative clustering as the complexity of divisive is higher than Agglomerative and has similar performance of Divisive approach.

Syllabus Topic : Expectation Maximization Clustering

6.2 Expectation Maximization Clustering

The Expectation Maximization (EM) algorithm can be applied to bring forth the best theory for the distributional parameters of some multi-modal data.

The best hypothesis for the distributional parameters is the maximum likelihood hypothesis – the one that maximizes the probability that this data we are looking at comes from K distributions, each with mean μ_k and variance σ_k^2 .

The Expectation Maximization (EM) Algorithm deal with missing labels by alternating between two steps :

1. **Expectation (E) :** Fix model and estimate missing labels
 2. **Maximization (M) :** Fix missing labels (or a distribution over the missing label(s)) and find the model that maximizes the expected log-likelihood of the data.
- Agglomerative Clustering**
- In Bottom up approach, every object is considered to be a cluster and in subsequent iterations they are merged into a single cluster. Therefore, it is also called as Hierarchical Agglomerative Clustering (HAC).

General EM Algorithm

The alternate steps until model parameters don't change much :

- o E step : Estimate distribution over labels given a certain fixed model.
- o M step : Choose new parameters for model to maximize expected log-likelihood of observed data and hidden variables.

Formal Setup for General EM Algorithm

Let $D = \{x^{(1)}, \dots, x^{(n)}\}$ be n observed data vectors.

Let $Z = \{z^{(1)}, \dots, z^{(n)}\}$ be n values of hidden variables (i.e. the cluster labels)

Log-likelihood of observed data given model :

$$L(\theta) = \log P(D|\theta) = \log \sum_Z P(D|Z|\theta)$$

Note : both θ and Z are unknown.

Where theta is a vector of unknown parameters.

In clustering, K (number of clusters) are not known initially. We can produce two results for the clustering :

- Hard clustering : In this method, data object or observation i belongs to only one cluster, so clusters can be produced like $\{C_1, C_2, \dots, C_K\}$.
- Soft clustering : The method says that data object or observation is more likely to belong to one of the K cluster by producing a probability distribution.

$$P[X_i \in C_k] = Y_{k,i} \quad \text{with } \sum_{k=1}^K Y_{k,i} = 1$$

From Fig. 6.2.1, we can observe that hard clustering is possible on the left side figure. But right side figure is not clear, so only based on probability samples can be classified to one of the class. This is the purpose of soft clustering. Algorithms like hierarchical clustering or K means produce some Hard clustering. The purpose of the EM clustering is to propose a Soft clustering method.

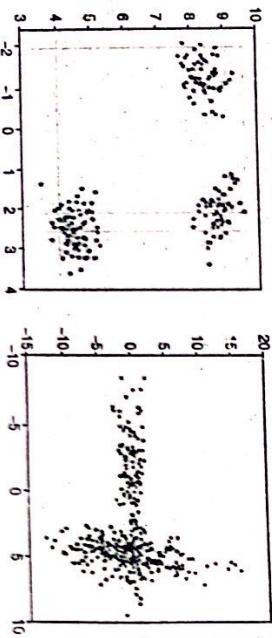


Fig. 6.2.1 : Visual point of view hard clustering and soft clustering

The basic steps for the algorithm are:

1. An initial guess is made for the model's parameter and a probability distribution is made. This is sometimes called the "E-Step" for the "Expected" distribution".
2. Newly observed data is fed into the model.
3. The chance distribution of the E-step is tweaked to let in the fresh information. This is sometimes called the "M-step".
4. Steps 2 through 4 are repeated until stability (i.e. a distribution that doesn't change from the E-step to the M-step) is reached.

Syllabus Topic : Agglomerative Clustering

6.3.1 Affinity

Affinity is a metric function of two arguments with the same dimensionality m.

The most common metrics supported by scikit-learn are:

1. Euclidean : Euclid stated that the shortest distance between two points is a line and thus it is predominantly known as Euclidean distance. It was often called Pythagorean metric since it is derived from the Pythagorean Theorem.

$$d_{\text{euclidean}} = \sqrt{\sum_{i=1}^d |P_i - Q_i|^2}$$

So the distance between two points is

$$\text{dist}(x, y), (a, b)) = \sqrt{(x-a)^2 + (y-b)^2}$$

2. Manhattan or City Block: The City block distance between two points, C and B, with d dimensions are calculated as

$$d_{\text{manhattan}} = \sum_{i=1}^d |P_i - Q_i|$$

3. The distance between two points measured along axes at right angles. In a plane with P_1 at (x_1, y_1) and P_2 at (x_2, y_2) , it is $|x_1 - x_2| + |y_1 - y_2|$.

The cosine distance is useful when we need a distance proportional to the angle between two vectors. If the direction is the same, the distance is null, while it is maximum when the angle is equal to 180° (meaning opposite directions).

$\cos(\theta) = \frac{\mathbf{P} \cdot \mathbf{Q}}{\ \mathbf{P}\ \ \mathbf{Q}\ }$	$i \cdot j = \sum_{k=1}^m i_k j_k$	$\ \mathbf{P}\ = \sqrt{\sum_{k=1}^m i_k^2}$	$\theta = \arccos \frac{\mathbf{P} \cdot \mathbf{Q}}{\ \mathbf{P}\ \ \mathbf{Q}\ }$
--	------------------------------------	--	--

Cosine Distance = 1 - Cosine Similarity

$$d_{\text{cosine}} (\tilde{x}_1, \tilde{x}_2) = 1 - \frac{\tilde{x}_1 \cdot \tilde{x}_2}{\|\tilde{x}_1\| \|\tilde{x}_2\|}$$

The cosine distance is useful when we need a distance proportional to the angle between two vectors. If the direction is the same, the distance is null, while it is maximum when the angle is equal to 180° (meaning opposite directions).

6.3.2 Strategy to Aggregate Different Clusters

Different approaches to defining the distance between clusters distinguish the different algorithms, but scikit learn supports the three most common ones i.e. Complete Linkage, Average Linkage and Ward's linkage.

1. Single-linkage

- Single Linkage clustering is also called as minimum method, the minimum distance from any object of one cluster to any object of another cluster is considered. In the single linkage method, $D(A,B)$ is computed as

$$D(A,B) = \min\{d(i,j)\}$$

Where object i is in cluster A and object j is in cluster B]

This measure of inter-group distance is illustrated in the Fig. 6.3.1.

Fig. 6.3.1 : Single-linkage clustering

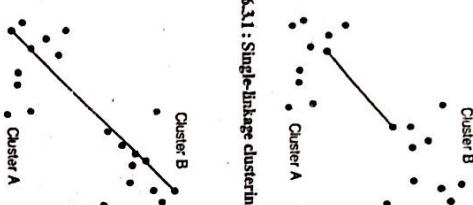


Fig. 6.3.2: Complete-linkage clustering

3. Average-linkage

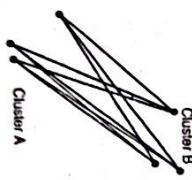
- In average-linkage clustering, we consider the distance between any two clusters A and B is taken to be the average of all distances between pairs of objects "i" in A and "j" in B, that is, the mean distance between elements of each cluster.

- In the average linkage method, $D(A,B)$ is computed as

$$D(A,B) = \text{Mean}\{d(i,j) : \text{Where object } i \text{ is in cluster A and object } j \text{ is in cluster B}\}$$

The measure is illustrated in the Fig. 6.3.2.

Fig. 6.3.3 : Average-linkage clustering



6.3.3 Dendograms

- An agglomerative clustering is typically visualized as a dendrogram as shown in Fig. 6.3.5 where each merge is represented by a horizontal line.
- Dendrogram is the aggregation of clusters from bottom to top.
- Initially leaf nodes are single clusters, which are merged till a single root node or cluster generated.

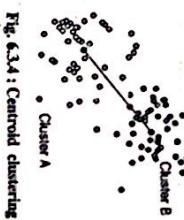


Fig. 6.3.5 : Dendrogram

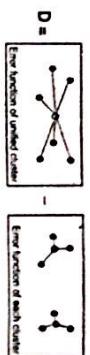
What is Dendrogram?

- Dendrogram : A tree data structure which illustrates hierarchical clustering techniques.

Each level shows clusters for that level.

- In this method, the error function is defined for every cluster. This error function is the average (RMSE) distance of each data point in a cluster to the centre of gravity in the cluster.
- A cluster at level l is the union of its children clusters at level $l+1$.

Syllabus Topic : Dendograms



- The distance (D) between two clusters is defined as the error function of the unified cluster minus the error functions of the individual clusters.
- The ward's linkage method computes the distance using formula,

$$\forall C_i, C_j, L_k = \sum_{x_i \in C_i} \sum_{x_j \in C_j} \|x_i - x_j\|^2$$

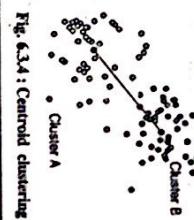
- The Ward's linkage supports only the Euclidean distance.

5. Centroid clustering

- In centroid method, the distance between two clusters is calculated by finding the distance between two centroids (i.e. mean value of cluster) of the clusters. At every step, two clusters are combined that have minimum centroid distance.
- In the centroid clustering method, $D(A,B)$ is computed as

$$D(A,B) = d(A_{\text{mean}}, B_{\text{mean}})$$

- Where A_{mean} is the mean value or centroid of cluster A and B_{mean} is the mean value or centroid of cluster B. The Fig. 6.3.4 illustrates centroid clustering.



A clustering of the data objects is obtained by cutting the dendrogram at the desired level, then each connected component form a cluster.

Unfortunately, scikit-learn doesn't support the dendrogram but SciPy provides some useful built-in functions.

The program in Python is given below to create dendrogram.

Program to visualize dendrogram

from sklearn.datasets import make_blobs

nb_samples = 25

X = make_blobs(n_samples=nb_samples, n_features=2, centers=3, cluster_std=1.5)

computing a distance matrix chosen a Euclidean metric

dmat = pdist(X, metric='euclidean')

linkage used is Ward

linkage = linkage(dmat)

```
from scipy.cluster.hierarchy import dendrogram
Y = linkage(X, method='ward')
plt.figure()
plt.title('Dendrogram')
from scipy.cluster.hierarchy import dendrogram
from scipy.cluster.hierarchy import linkage
plt.title('Dendrogram')
plt.show()
```

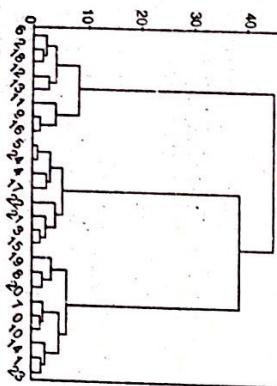


Fig. 6.3.6 : Snapshot of resulting plot of dendrogram

6.3.4 Agglomerative Clustering In Scikit-learn

Consider a more complex dummy dataset with 4 centers.

from sklearn.datasets import make_blobs

from sklearn.cluster import AgglomerativeClustering

import numpy as np

We have now a data set with 8000 random samples around the 4 center points with a standard deviation of 2 for np.random.seed(42)

```
nb_samples = 3000
X = make_blobs(n_samples=nb_samples, n_features=2, centers=4, cluster_std=2.0)
```

visualize the data set

plt.scatter(X[:, 0], X[:, 1])

plt.show()

A graphical representation is shown in the following Fig. 6.3.7.



Fig. 6.3.7 : snapshot of resulting plot of data representation

X is the dataset, and y is the label of each data point, according to the sample generation. Agglomerative Clustering uses the method fit_predict() to train the model and transform the original dataset. The number of parameters is set to 4 using the n_clusters parameter, while the affinity is set to "Euclidean" (distance between the data points). Finally linkage parameter is set to 'complete'.

```
mc = AgglomerativeClustering(n_clusters=4, affinity='euclidean', linkage='complete')
y = mc.fit_predict(X)
ll.scatter(X[:, 0], X[:, 1], c=y, cmap='rainbow')
```

A graphical representation is shown in the following Fig. 6.3.8.

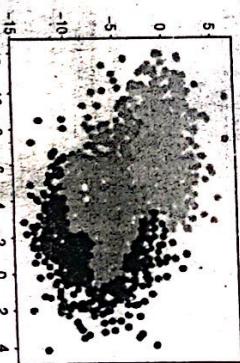


Fig. 6.3.8 Snapshot of resulting plot of complete linkage

Linkage parameter is set to "average".

mc = AgglomerativeClustering(n_clusters=4, affinity='euclidean', linkage='average')

y = mc.fit_predict(X)

ll.scatter(X[:, 0], X[:, 1], c=y, cmap='rainbow')

The resultant plot is shown in Fig. 6.3.9 :



Fig. 6.3.9 : Snapshot of resulting plot of average linkage

- The last method, which is often the best (it's the default one), is Ward's linkage, that can be used only with a Euclidean metric (also the default one):

```
# Now we can try to impose a constraint with different values for k
from sklearn.cluster import KMeans
from sklearn.neighbors import NeighborsGraph
n_neighbors = 20
graph = NeighborsGraph(n_neighbors=n_neighbors)
X = graph.fit_transform(X)
kmeans = KMeans(n_clusters=2, random_state=42, n_init='auto')
y_kmeans = kmeans.fit_predict(X)

# Let's plot the result
plt.scatter(X[:, 0], X[:, 1], c=y_kmeans, cmap='rainbow')
```

The resultant plot is given in Fig. 6.3.10.

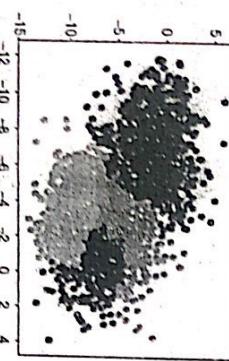


Fig. 6.3.10 : Snapshot of resulting plot of ward linkage

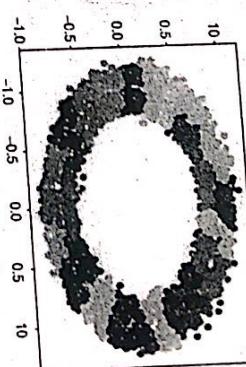
Syllabus Topic : Connectivity Constraints

6.3.5 Connectivity Constraints

- An interesting aspect of Agglomerative Clustering is that *connectivity constraints* which can be added to this algorithm where only adjacent clusters can be merged together and those clusters which are distant i.e. non-adjacent are skipped. Connectivity matrix can be used as a constraint by SciKit-learn and find the clusters to merge. These constraints are useful to impose a certain local structure, but they also make the algorithm faster, especially when the number of the samples is high.
- K-Nearest Neighbour graph function is a common method used.
- Consider circular dummy dataset in the example given below.

```
# Now we can try to impose a constraint with different values for k
from sklearn.cluster import KMeans
from sklearn.neighbors import NeighborsGraph
n_neighbors = 20
graph = NeighborsGraph(n_neighbors=n_neighbors)
X = graph.fit_transform(X)
kmeans = KMeans(n_clusters=2, random_state=42, n_init='auto')
y_kmeans = kmeans.fit_predict(X)

# Let's plot the result
plt.scatter(X[:, 0], X[:, 1], c=y_kmeans, cmap='rainbow')
```



```
# Now we can try to impose a constraint with different values for k
from sklearn.cluster import KMeans
from sklearn.neighbors import NeighborsGraph
n_neighbors = 20
graph = NeighborsGraph(n_neighbors=n_neighbors)
X = graph.fit_transform(X)
kmeans = KMeans(n_clusters=2, random_state=42, n_init='auto')
y_kmeans = kmeans.fit_predict(X)

# Let's plot the result
plt.scatter(X[:, 0], X[:, 1], c=y_kmeans, cmap='rainbow')

# Imposing a constraint based on K-Nearest Neighbours, allows controlling how the agglomeration creates new clusters
# and can be a powerful tool for tuning the models, or for avoiding elements whose distance is large in the original space
# and can be taken into account during the merging phase. It is helpful in clustering the images.
```

Syllabus Topic : Introduction to Recommendation Systems - Native User based Systems

6.4 Introduction to Recommendation Systems

- Recommendation system is used to predict future preferences of a set of items and also recommend the top n item users.
- Many e-commerce websites uses recommendation system to increase their sales by using the power of data.
- This system finds the user's interest to recommend items.

Traditionally, there are two methods to construct a recommendation system :

1. User or Content-based recommendation
2. Collaborative Filtering

1. **Content-based recommendation system:** It makes recommendations using a user's item and profile features. To find the user's future interest, the system needs the information about user's past interest in the items. So only similar items can be grouped together based on their features.
2. **Collaborative Filtering systems:** It filters the items in which we are interested. This system is based on the assumption that if user1 likes item X and user2 likes the same item X plus another item Y, then user1 may be interested in item Y. So system can use historical data to predict new interactions.

6.4.1 Naïve User based Systems

Let us consider a set of users represented by feature vectors :

$$U = \{\bar{u}_1, \bar{u}_2, \dots, \bar{u}_n\} \text{ where } \bar{u}_i \in \mathbb{R}^m$$

Assume the features are age, gender, interests, and so on. Set of item I is :

$$I = \{i_1, i_2, \dots, i_n\}$$

Every user is associated with a subset of items or items can be viewed or feedback has been taken for those items. So there is a relationship between user and items.

$$g(\bar{u}) \rightarrow \{i_1, i_2, \dots, i_k\} \text{ where } k \in (0, m)$$

In a user-based system, the users are periodically clustered (normally using a K-Nearest Neighbour approach), and therefore, considering a generic user \bar{u} (also new), we can immediately determine all the users who are similar (before neighbours) to our sample.

$$B_i(\bar{u}) = \{\bar{u}_1, \bar{u}_2, \dots, \bar{u}_k\}$$

At this point, we can create the set of suggested items using the relation previously introduced :

$$I_{suggested}(\bar{u}) = \left\{ \begin{array}{l} i_j \\ \text{if } g(\bar{u}_j) \text{ where } \bar{u}_j \in B_i(\bar{u}) \end{array} \right\}$$

This set contains all the unique products positively rated or bought by the neighbourhood.

Syllabus Topic : Content Based Systems

6.4.2 Content Based Systems

Recommendations are based on information on the content of items rather than on other users' opinions.

Uses machine learning algorithms to induce a profile of the user's preferences from examples based on content features.

- o No need for data about other users.
- o No cold-start or sparsity problems.
- o Able to recommend to users with unique tastes.
- o No first-rater problem.
- **Cold Start :** enough users in the system to find a match.
- **Sparsity :** The user/ratings matrix is sparse, and it is hard to find users that have rated the same items.

First Rater : Not for an item that has not been previously rated.
This is one of the simplest methods which is based on products and modelled as feature vectors :

$$I = \{\bar{i}_1, \bar{i}_2, \dots, \bar{i}_n\} \text{ where } \bar{i}_j \in \mathbb{R}^m$$

The features can also be categorized like users. For example, the type of a book or a movie, and they can be used together with numerical values (like price, length, number of positive reviews, and so on) after encoding them.

Syllabus Topic : Model Free Collaborative Filtering

6.4.3 Model Free Collaborative Filtering

It is based on rating among the users and "lightweight" filtering approach. It is computed in-memory without the use of complex algorithms like ALS (Alternating Least Squares) which executes in parallel environment.

Assume that we have M products and N users, then user-preference sparse matrix is given below.

$$U_{prod} = \begin{pmatrix} & & & & \\ u_0^{(0)} & \cdots & u_0^{(M)} & & \\ & \ddots & \ddots & \ddots & \\ & & \vdots & & \\ u_N^{(0)} & \cdots & u_N^{(M)} & & \end{pmatrix}$$

Each element from the above represents the rating given by the user i to the item j .

0 means there is no rating.

Two kinds of ratings or feedbacks

1. Explicit feedback : An actual rating (like 3 stars or 8 out of 10)
2. Implicit feedback : A page view, song play, movie play and so forth

After the user preference matrix, compute the affinity matrix based on a similarity metric using Euclidean as given below :

$$A = \begin{pmatrix} d(u_0, u_0) & \cdots & d(u_0, u_M) \\ \vdots & \ddots & \vdots \\ d(u_N, u_0) & \cdots & d(u_N, u_M) \end{pmatrix}$$

For explicit feedback, the Euclidean metric is preferable, but when it's implicit (like page views, song/movie plays, etc.), it has a different meaning and need another metric function. For example, if we have a binary encoding like 0 and 1 (0 means no interaction and 1 means single or multiple interaction), then Hamming distance is more appropriate. Our goal is to cluster the users according to their rating. After clustering, we can check which products have higher rating for a given user and therefore are more likely to be bought.

Once the top neighbours determined, union of most rated neighbour products is used to suggest the list of top neighbours :

$$P_{neig}(i) = \bigcup_{j \neq i} P(i)$$

The code to create a distance matrix using sklearn for 10 users and 10 items.

```
from scipy import sparse
from sklearn.metrics.pairwise import pairwise_distances
import numpy as np
```

```
# random seed for reproducibility
```

```
np.random.seed(42)
```

```
user_id = 10
```

```
# create dummy user-item dataset
```

```
nb_users = 10
```

```
nb_products = 10
```

```
user_rating = 5
```

```
user_cited_products = 8
```

```
product = np.zeros(nb_products, dtype=np.int) + user_rating
```

```
preferences = dict((i, user_cited_products * product) for i in range(nb_products))
```

```
# regular preference matrix
```

```
matrix = np.array([preferences[i] for i in range(nb_products)])
```

```
# random user-item distances
```

```
# Create n random products
```

```
# Products = np.zeros(nb_products, dtype=np.int) + nb_products * sum(np.array([0, 1]))
```

```
product = np.zeros(nb_products, dtype=np.int) + nb_products * sum(np.array([0, 1]))
```

```
# Regular preference matrix
```

```
matrix = np.array([product[i] for i in range(nb_products)])
```

```
# Compute pairwise distances
```

```
distance_matrix = pairwise_distances(matrix, metric='euclidean')
```

```
# Sort distances
```

```
sorted_distances = np.argsort(distance_matrix, axis=-1)
```

```
# Compute distances
```

```
[0 1 2 3 4 5 6 7 8 9]
```

```
[0 1 2 3 4 5 6 7 8 9]
```

```
[0 1 2 3 4 5 6 7 8 9]
```

```
[0 1 2 3 4 5 6 7 8 9]
```

```
[0 1 2 3 4 5 6 7 8 9]
```

Output

Print the top-5 similar users or 5 nearest distances[[test user]]:[1][0:5]

```
print(d[0:5])
```

Output

```
[1 2 3 4 5]
```

Syllabus Topic : Singular Value Decomposition

6.4.4 Singular Value Decomposition

Singular value decomposition (SVD) is a method of decomposing a matrix into three other matrices.

Given any $n \times n$ matrix A, algorithm to find matrices U, V, and W such that

$$A = U W V^T$$

U is $m \times n$ and orthonormal

W is $n \times n$ and orthogonal

$$(A) = (U) \begin{pmatrix} w_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & w_n \end{pmatrix} (V)^T$$

By assuming $m \geq n$ rewrite above equation using summation notation

$$a_{ij} = \sum_{k=1}^n u_{ik} s_k v_{jk}$$

Simplify the expression into a single summation. The variables, s_k , are called singular values and are normally arranged from largest to smallest.

$$s_{i+1} \leq s_i$$

The columns of U are called left singular vectors, while those of V are called right singular vectors.

We know that U and V are orthogonal, that is:

$$U^T U = V V^T = I$$

Where I is the identity matrix (the diagonals of the identity matrix are 1, other values being 0).

svd() function calculates SVD.

svd() function takes matrix as input and returns U, Sigma and V^T

Output

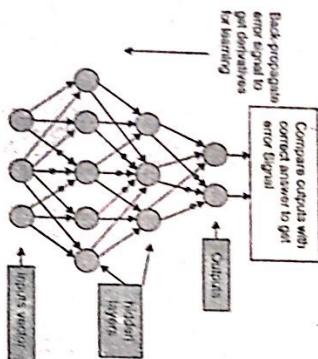
Output

6.5.2(A) Multilayer Perceptron with Back-propagation : First Deep Learning Model

(Rumelhart, Hinton, Williams 1986)

Multilayer Perceptrons with Back-propagation

First deep learning model (Rumelhart, Hinton, Williams 1986)

**or Drawbacks of Back-propagation based Deep Neural Networks**

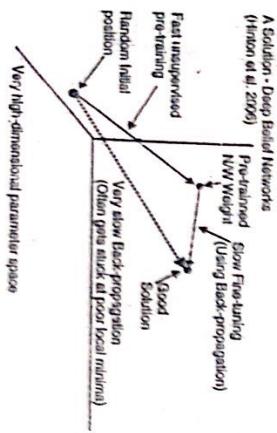
- They are discriminative models.
- All the information is collected from labels.
- But the labels are not able to give all information.
- It requires a large amount of labeled data.
- Gradient descent with random initialization leads to poor local minima.

6.5.2(B) Deep Belief Networks (Hinton et al. 2006)

Pre-training of network is required before using back-propagation.

Use the weights of the pre-trained network as the initial point for the traditional back-propagation.

- o This leads to quicker convergence to a better solution.
- o The training is fast; fine-tuning can be slow.

**6.5.2(C) Five Popular and Widely-used Deep Learning Architectures**

1. Convolutional Neural Networks
2. Recurrent Neural Networks
3. Autoencoders
4. Generative Adversarial Networks
5. RNNs

Syllabus Topic : Building Blocks of Deep Networks**6.5.2(D) Building Blocks of Deep Networks**

- Activation functions : linear, ReLU, sigmoid, tanh, etc.
- Layers : Fully connected, convolutional and pooling, recurrent (Problem: The graph is not acyclic. To compute the gradients, we unroll the computational cycle over timesteps, and backpropagate through this structure), recurrent (skip connections over layers), etc.

Review Questions

- Q.1 Explain Hierarchical Clustering basic algorithm.
- Q.2 Explain Expectation Maximization Clustering algorithms in detail.
- Q.3 What are various strategy to Aggregate Different Clusters?
- Q.4 Explain Model Free Collaborative Filtering.
- Q.5 Write short note on Recommendation Systems
- Q.6 What is Deep Learning? Explain.

